

UNIVERSIDAD AUTÓNOMA DE MADRID  
ESCUELA POLITÉCNICA SUPERIOR



TRABAJO FIN DE MÁSTER

# Algoritmos Probabilísticos Automáticos para Predicción a partir de Señales de Centrales Nucleares

Máster Universitario en Ingeniería de Telecomunicación

Autor:Pablo Ramírez Hereza  
Tutor:Daniel Ramos Castro

FECHA:SEPTIEMBRE 2019



# ALGORITMOS PROBABILÍSTICOS AUTOMÁTICOS PARA PREDICCIÓN A PARTIR DE SEÑALES DE CENTRALES NUCLEARES

AUTOR: Pablo Ramírez Hereza  
DIRECTOR: Daniel Ramos Castro

AUDIAS  
Dpto. de Tecnología electrónica y de las Comunicaciones  
Escuela Politécnica Superior  
Universidad Autónoma de Madrid  
SEPTIEMBRE 2019



## Resumen

Este Trabajo de Fin de Máster surge con la colaboración del grupo *AUDIAS*, de la Escuela Politécnica Superior (UAM), y una empresa gestora de centrales nucleares. Así pues, en el inicio del proyecto, varios objetivos son definidos por la empresa, en concreto, la adquisición de un mayor conocimiento de las relaciones entre los principales elementos químicos que actúan en un reactor nuclear, así como la predicción, a lo largo del ciclo de actuación del reactor, de los valores de los elementos más relevantes en la formación del residuo final en el interior de dicho reactor. De esta forma, en un futuro, la empresa pretende poder predecir dicho residuo y definir protocolos de actuación en función de dicha predicción.

Para ello, el proyecto se ha dividido en varias etapas. Primero se ha realizado un análisis de las señales, de naturaleza química y procedentes de sensores, involucradas en el proceso de generación de energía, incluyendo además de un análisis de su frecuencia de adquisición y de la correlación entre ellas. Posteriormente, se ha entrenado una Red Bayesiana para la predicción de las señales de interés. Tras la primera aproximación, se ha realizado un proceso, en colaboración con la empresa, de mejora del modelo y de transformación de los datos mediante técnicas ampliamente utilizadas en procesamiento de señales, como son la gaussianización e interpolación. Por último, se ha propuesto un modelo más complejo, basado en una aproximación dinámica de la Red Bayesiana previamente implementada, que explotará información temporal de las variables.

Una vez finalizadas las aproximaciones propuestas, se ha implementado una primera versión de una interfaz de usuario que permite realizar predicciones, para cualquier central, a partir de la Red Bayesiana, ya entrenada con los datos disponibles o con nuevos datos introducidos en la aplicación.

De esta forma, este trabajo supone una transferencia de conocimiento y de tecnología a la empresa, obteniendo unos resultados en la predicción que cumplen las expectativas iniciales del proyecto.

## Palabras Clave

Redes Bayesianas, Redes Bayesianas Dinámicas, gaussianización, predicción, reactor nuclear, química nuclear.



# Abstract

## Abstract

This Final Master in Science Thesis arises with the collaboration between the group *AUDIAS*, of the “Escuela Politécnica Superior”(UAM), and a nuclear power plant management company. In this way, some objectives have been defined by the company at the beginning of the project. Specifically, the acquisition of a larger knowledge about the relations between the main chemical elements that act in a nuclear reactor, and the prediction of the values of the most relevant elements in the generation of the final residue inside the reactor, to allow, in the future, the prediction of the final residue and the definition of actions according to these predictions.

In order to reach these objectives, the project has been divided in different steps. First of all, an analysis of the signals, of chemical nature and captured by sensors, involved in the process of energy generation, its acquisition frequency and its correlation, has been done. Later, a Bayesian Network has been trained for the prediction of the signals of interest for each nuclear power plant. After this first approximation, a process of model improvement and data transformation has been done, using very popular techniques of signal processing as Gaussianization and interpolation. Finally, a more complex model is proposed. This model is based on a dynamic approximation of the bayesian network already implemented which will exploit temporal information of these variables.

Once both models have been trained and evaluated, a first version of an user interface has been implemented. This interface allows to make predictions, for every nuclear power plant, using the Bayesian network, trained with the available data or with fresh data to be loaded in the application.

Then, this work involves a knowledge and technology transfer to the company. The results obtained in the prediction reach the initial expectations of the project.

## Key Words

Bayesian Networks, Dinamyc Bayesian Networks, gaussianization, prediction, nuclear reactor, nuclear chemistry.





# Agradecimientos

*Querido abuelo,*

*Sé que si estuvieses hoy aquí, leerías estos agradecimientos a cada persona que te visitara, que seguirías sin entender la necesidad de realizar un máster para ser ingeniero, y que serías de las personas que mejor reconocería el esfuerzo y dedicación necesarios para finalizar una etapa que, sin lugar a duda, ha sido una de las más complicadas de mi vida.*

*Por aquí, algunas cosas van cambiando. Papá y mamá se van a Valencia, espero que por fin a descansar del trabajo que le siguen dando todavía sus hijos. Les voy a necesitar y a echar de menos más de lo que ellos creen, pero sé con certeza que están delante de una etapa muy enriquecedora que ya se merecen. Así que Guillermo, Fatimita y yo nos quedamos aquí, y tranquilo, que aunque vaya un poco a mi bola, ellos saben que siempre estaré a su lado cuando lo necesiten. Al igual que ellos siempre han hecho conmigo. Además, aunque siempre echaremos de menos bajar a visitarte, seguimos teniendo la compañía del tío Pipo y de su nuevo apartamento.*

*Por otro lado, hay cosas que no cambian, el grupo de amigos sigue muy unido y, aunque cada vez sea más complicado reunirse todos y que, para que engañarnos, cada vez damos menos guerra en el campo de fútbol, sé que siempre puedo contar con ellos cuando necesito desconectar del mundo. Tania sigue siendo el mayor apoyo que tengo, y, a medida que cerramos etapas juntos, más cerca estamos en conseguir todos aquellos retos de los que hablamos y que en otro momento te contaré. Sigues muy presente en nuestra relación, e incluso, aunque no te llegaron a conocer, sigues presente en conversaciones con sus padres Guillermo y Maria Luisa. Ojala les hubieses conocido. Nunca podré agradecerles lo suficiente, a ellos y también a Andrés, Diana y Manuelita el acogerme como uno más de la familia.*

*Todos, incluyéndote a ti, sin estar aquí, de una forma u otra, han sido los pilares fundamentales para la finalización de esta etapa y, por lo tanto de este trabajo. GRACIAS.*

*Te quiero,  
Pablo.*

*Juan Hereza Domínguez.*



*Mi amigo, mi profesor, mi modelo a seguir. Mi abuelo.*



# Índice general

<b>Índice de figuras</b>	<b>XI</b>
<b>Índice de cuadros</b>	<b>XIII</b>
<b>1. Introducción</b>	<b>1</b>
1.1. Motivación del proyecto . . . . .	1
1.2. Objetivos y enfoque . . . . .	2
1.3. Metodología y plan de trabajo . . . . .	3
1.4. Organización de la memoria . . . . .	3
<b>2. Estado del Arte</b>	<b>5</b>
2.1. Teoría de la probabilidad . . . . .	5
2.2. Independencia marginal e independencia condicional . . . . .	6
2.3. Modelos Probabilísticos Gráficos. Redes Bayesianas . . . . .	7
2.3.1. Introducción . . . . .	7
2.3.2. Redes Bayesianas . . . . .	7
2.3.3. Inferencia en Redes Bayesianas . . . . .	8
2.3.4. Aprendizaje en Redes Bayesianas . . . . .	10
2.4. Introducción a las Redes Bayesianas Dinámicas . . . . .	12
2.5. Gaussianización de datos . . . . .	13
2.5.1. Gaussianización basada en ecualización de histogramas . . . . .	13
2.5.2. Feature Warping . . . . .	13
<b>3. Diseño</b>	<b>15</b>
3.1. Base de datos . . . . .	15
3.2. Diseño del sistema . . . . .	16
3.3. Diseño de la evaluación . . . . .	18
3.4. Diseño de la Aplicación . . . . .	18
<b>4. Desarrollo del Proyecto</b>	<b>21</b>
4.1. Analisis Preliminar . . . . .	21

4.1.1. Análisis del comportamiento de las señales . . . . .	21
4.1.2. Análisis comparativo temporal . . . . .	25
4.1.3. Análisis de la correlación mediante Scatter Plots . . . . .	26
4.2. Primera aproximación . . . . .	28
4.3. Modificación del modelo y de los datos . . . . .	29
4.3.1. Modificaciones directamente sugeridas por <i>La Empresa</i> . . . . .	29
4.3.2. Gaussianización de los datos . . . . .	30
4.3.3. Interpolación de datos . . . . .	31
4.4. Aproximación Dinámica . . . . .	32
4.5. Desarrollo de la Aplicación . . . . .	32
<b>5. Experimentos y Resultados</b>	<b>35</b>
5.1. Resultados Red Bayesiana inicial . . . . .	35
5.2. Resultados Red Bayesiana Final . . . . .	38
5.3. Resultados Red Bayesiana Dinámica . . . . .	43
<b>6. Conclusiones y trabajo futuro</b>	<b>49</b>

# Índice de figuras

2.1. Ejemplo de Red Bayesiana . . . . .	8
2.2. Ejemplo de Red Bayesiana Dinámica. . . . .	12
2.3. Función Probit, o inversa de la normal, para gaussianización de datos con distribución uniforme . . . . .	13
2.4. Proceso de Gaussianización mediante Feature Warping propuesto en [13]. Esta imagen ha sido extraída del mismo documento . . . . .	14
3.1. Red Bayesiana Final . . . . .	17
3.2. Red Bayesiana Dinámica Final . . . . .	17
3.3. Ejemplo Hugin sin evidencias de la Red Bayesiana inicial . . . . .	19
3.4. Ejemplo Hugin con evidencias de la Red Bayesiana inicial . . . . .	19
4.1. Señal <i>Control 1</i> . . . . .	22
4.2. Señal <i>Medida 1</i> . . . . .	24
4.3. Análisis comparativo temporal para la central <i>Planta 1</i> . . . . .	26
4.4. Ejemplo del análisis realizado. La gráfica de arriba representa un diagrama de dispersión. La gráfica de abajo representa el coeficiente a corto plazo, en azul si $p\text{-valor} < 0.05$ , mediante una cruz roja en el caso contrario. . . . .	27
4.5. Red Bayesiana Inicial . . . . .	28
4.6. Desarrollo del modelo inicial . . . . .	29
4.7. Limpieza de <i>Control 1</i> para <i>Planta 1</i> el ciclo 24. (a) Señal antes del proceso. (b) Señal después del proceso. . . . .	30
4.8. Derivación de <i>Control 2</i> en <i>Planta 1</i> el ciclo 22. (a) Señal inicial. (b) Derivada . . . . .	30
4.9. Proceso de Gaussianización en señal <i>Medida 1</i> , <i>Planta 3</i> ciclo 24 . . . . .	31
4.10. Diagrama final del proceso de desarrollo . . . . .	32
4.11. Aplicación AUDIAS para el proyecto . . . . .	33
5.1. Predicción con modelo inicial para <i>Planta 3</i> Ciclo 26 . . . . .	36
5.2. Predicción con modelo inicial para <i>Planta 2</i> Ciclo 19 . . . . .	36
5.3. Predicción con modelo inicial para <i>Planta 4</i> Ciclo 25 . . . . .	37
5.4. Predicción final para <i>Planta 1</i> Ciclo 24 . . . . .	39
5.5. Predicción final para <i>Planta 3</i> Ciclo 25 . . . . .	39

5.6. Predicción final para <i>Planta 5</i> Ciclo 20 . . . . .	40
5.7. Comparativa RMSE medio tras las modificaciones propuestas . . . . .	42
5.8. RMSE medio para todas las centrales en función de N . . . . .	43
5.9. Predicción con DBN en <i>Planta 2</i> Ciclo 23 . . . . .	44
5.10. Predicción en <i>Planta 5</i> Ciclo 20 . . . . .	44
5.11. Comparativa RMSE medio para todas las centrales para todas las técnicas . . . .	47

## Índice de cuadros

5.1. Resultados obtenidos primera aproximación . . . . .	38
5.2. Resumen resultados Finales <i>Medida 1</i> . . . . .	45
5.3. Resumen resultados Finales <i>Medida 2</i> . . . . .	45
5.4. Resumen resultados Finales <i>Medida 3</i> . . . . .	45





# 1

## Introducción

### 1.1. Motivación del proyecto

---

Este Trabajo de Tin de Máster se enmarca en la predicción de variables de interés a partir de señales temporales procedentes de centrales nucleares en el contexto de la colaboración entre el grupo AUDIAS de la Escuela Politécnica Superior, de la Universidad Autónoma de Madrid, y una empresa dedicada a la gestión industrial de centrales nucleares. Dicha empresa será denominada como *La Empresa* a lo largo de este documento.

Entre las actividades de *La Empresa*, se encuentra la generación de protocolos y recomendaciones sobre el procedimiento de generación de electricidad mediante la obtención de vapor de agua a partir de reacciones nucleares de fisión, lo que es llamado el circuito primario del reactor nuclear, en este documento referido como el *primario*. Dicho proceso es complejo y requiere de múltiples elementos físicos y químicos a tener en cuenta. El funcionamiento del *primario* consiste en la repetición de *ciclos* temporales en los que se produce la carga de combustible, el funcionamiento del reactor y, por último, la limpieza del mismo. Así pues, dentro de la química involucrada en el proceso, existen ciertos aspectos en los cuales *La Empresa* tiene interés en profundizar.

Esta entidad, pese a cumplir escrupulosamente con la normativa aplicable y los estándares de calidad, pretende minimizar la liberación de productos de corrosión mediante el correcto uso y monitorización de determinados aditivos químicos. De esta forma, a lo largo de cada ciclo de funcionamiento, diferenciaremos las señales químicas en señales *Control* y señales *Medida*.

Las señales, o variables *Control*, representan los valores de los aditivos, o medidas realizadas sobre ellos, y los parámetros controlables durante el funcionamiento del reactor. Dichas señales pueden estar relacionados químicamente entre sí e interaccionar en consecuencia, no estando completamente claro para *La Empresa* el impacto de éstos sobre las señales, o variables *Medida*, subproductos del proceso completo directamente relacionados con la corrosión final, y cuyo valor se pretende predecir en función de dichas señales, o variables, de control.

Los resultados esperados de este Trabajo de Fin de Máster son, por tanto, el análisis de datos de la química del *primario* y la implementación de algoritmos de predicción con el objetivo de obtener información detallada de las relaciones entre series temporales de interés para *La Empresa* a lo largo del proceso. Para ello, se plantea el uso de algoritmos basados en modelos gráficos probabilísticos, y en particular Redes Bayesianas, que permitan establecer la relación entre señales de forma gráfica, visualmente sencilla e interpretable por *La Empresa*. En concreto, se utilizarán Redes Bayesianas Gaussianas debido a sus propiedades analíticas. Adicionalmente, se hará uso de aproximaciones dinámicas de las Redes Bayesianas utilizadas, para añadir información temporal al modelo gráfico final.

De esta forma, el impacto científico de este proyecto es alto pues, hasta donde alcanza nuestro conocimiento, no existe ninguna aplicación de métodos probabilísticos gráficos de aprendizaje automático que utilicen datos de química del *primario* en centrales nucleares. Por otro lado, se espera que los resultados de este trabajo sean útiles para la mejora de los procedimientos de *La Empresa*, en relación con el control de la calidad y el aumento de la eficiencia energética en centrales nucleares, lo cual implicará una alta innovación tecnológica en el ámbito del análisis estadístico y de la ciencia de datos en el contexto de esta aplicación.

## 1.2. Objetivos y enfoque

---

El objetivo principal de este trabajo es, por lo tanto, el desarrollo para cada central de un sistema basado en métodos probabilísticos de aprendizaje automático y en técnicas de procesamiento de señal, que permita la representación de las dependencias entre las diferentes señales químicas presentes en el *primario* y una predicción de las variables *Medida*.

Para esto, se han definido una serie de objetivos parciales:

1. Estudio de los datos proporcionados por *La Empresa*.
2. Estudio, implementación y entrenamiento de una Red Bayesiana que relacione las señales químicas y permita la predicción de las variables de interés.
3. Aplicación de técnicas de procesamiento de señal para una mejor representación de los datos por el modelo gráfico y así mejorar el rendimiento en la predicción.
4. Uso de aproximaciones dinámicas de la Red Bayesiana.
5. Implementación de una interfaz gráfica de usuario, utilizable por la empresa que le permita el uso del modelo probabilístico propuesto.

### 1.3. Metodología y plan de trabajo

---

Para el correcto cumplimiento de los objetivos propuestos, se sigue un plan de trabajo definido por las siguientes etapas:

1. Etapa meramente organizativa, basada en la planificación del proyecto. En esta etapa se fijan los plazos (*deadlines*) e hitos (*milestones*) del proyecto, y se determinan los procedimientos para su seguimiento.
2. Análisis exhaustivo de las señales químicas que confoman la base de datos.
3. Estudio e implementación de modelos gráficos probabilísticos, en concreto Redes Bayesianas Gaussianas. De esta forma, esta etapa nos aportará una representación de las dependencias de las señales involucradas en el *primario*, además de los primeros resultados de predicción de las señales medida.
4. Estudio e implementación de técnicas de gaussianización, en concreto basadas en normalización y ajuste de histogramas.
5. Estudio e implementación de una aproximación dinámica de la Red Bayesiana anterior, para permitir la consideración de dependencias a lo largo del tiempo.
6. Implementación de la aplicación entregable para *La Empresa*.
7. Comparativa entre los resultados obtenidos para diferentes centrales nucleares.
8. Redacción de este documento y defensa del trabajo expuesto.

### 1.4. Organización de la memoria

---

Esta memoria ha sido dividida de tal forma que el lector siga el flujo de trabajo del proyecto, entendiendo los conceptos clave en los que éste se basa. De esta forma, este documento se divide en varios capítulos. El primero, ofrece una introducción del proyecto y de su contexto. En el segundo se ofrecen las bases teóricas necesarias para el entendimiento de la solución propuesta en el trabajo, ofreciendo una explicación de conceptos básicos de probabilidad, de Redes Bayesianas y de técnicas de gaussianización. Por otro lado, la tercera sección describe el diseño final de los modelos propuestos, mientras que el capítulo 4 ofrece una descripción del proceso de desarrollo, describiendo las principales consideraciones para la obtención de los modelos finales. Por último, se presentan los principales resultados obtenidos, las conclusiones y el trabajo futuro, en los últimos dos capítulos de este documento.



# 2

## Estado del Arte

### 2.1. Teoría de la probabilidad

---

En esta sección veremos un breve repaso de las reglas fundamentales de la teoría de la probabilidad [1][2], base de los modelos probabilísticos presentados en el trabajo.

Dadas las variables aleatorias  $A = \{a_1, a_2, \dots, a_N\}$  y  $B = \{b_1, b_2, \dots, b_M\}$ :

1. Definimos la **regla de la suma** como:

$$P(X) + P(\bar{X}) = 1 \quad (2.1)$$

2. Definimos la **probabilidad conjunta** de ambos eventos mediante la que es denominada la **regla del producto** como:

$$P(A, B) = P(A|B)P(B). \quad (2.2)$$

Dadas las dos reglas fundamentales de la teoría de la probabilidad podemos definir:

1. **Probabilidad marginal** de A, en ocasiones presentada como la regla de la suma:

$$P(A) = \sum_b P(A, B) = \sum_i^M P(A|B = b_i)P(B = b_i) \quad (2.3)$$

2. **Teorema de Bayes**: A partir de la regla del producto y de la propiedad de simetría  $P(X, Y) = P(Y, X)$  obtenemos directamente la relación:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (2.4)$$

3. **Regla de la cadena**: Podemos definir la probabilidad conjunta de múltiples variables como:

$$P(X_1, X_2, \dots, X_R) = P(X_1)P(X_2|X_1)P(X_3|X_2, X_1) \dots P(X_R|X_1, X_2, \dots, X_3) \quad (2.5)$$

## 2.2. Independencia marginal e independencia condicional

---

Dadas dos variables aleatorias  $A = \{a_1, a_2, \dots, a_N\}$  y  $B = \{b_1, b_2, \dots, b_M\}$ :

Se dice que  $A$  es independiente marginalmente de  $B$  si se cumple:

$$P(A|B) = P(A) \quad (2.6)$$

Sustituyendo en (2.2), obtenemos:

$$P(A, B) = P(A)P(B) \quad (2.7)$$

Dada una tercera variable aleatoria  $Z = \{z_1, z_2, \dots, z_O\}$ . Se dice que  $A$  es independiente de  $B$  conociendo  $Z$ , es decir condicionalmente independiente, si se cumple que:

$$P(A|B, Z) = P(A|Z) \quad (2.8)$$

Sustituyendo en (2.2), obtenemos:

$$P(A, B|Z) = P(A|Z)P(B|Z) \quad (2.9)$$

## 2.3. Modelos Probabilísticos Gráficos. Redes Bayesianas

---

### 2.3.1. Introducción

Se ha visto en el capítulo anterior las principales reglas que forman la teoría de la probabilidad: la regla de la suma y la regla del producto.

Todos los modelos probabilísticos, su aprendizaje e inferencia se basan en la aplicación de estas dos ecuaciones. Pero, pese a que se podrían resolver problemas complejos probabilísticos analíticamente, es muy útil el uso de representaciones mediante diagramas de las distribuciones de probabilidad, pues, tal y como se define en [1] [3]:

1. Proveen una forma sencilla de visualizar la estructura de un modelo probabilístico.
2. Las propiedades del modelo, incluyendo las propiedades de independencia condicional, pueden ser obtenidas directamente del grafo.
3. Los cálculos complejos necesarios para la inferencia y el aprendizaje pueden expresarse como manipulaciones gráficas.
4. Un grafo específico establece las relaciones entre variables independientemente de las distribuciones de dichas variables.

Un modelo gráfico está compuesto por *nodos*, que representan las variables aleatorias existentes en el modelo (cada una con una distribución de probabilidad propia, discreta o continua), y *enlaces* que expresan las relaciones probabilísticas entre dichas variables. En función de dichos enlaces encontramos dos principales grupos de grafos probabilísticos:

1. **Modelos gráficos dirigidos:** Son los grafos en los cuales los enlaces se encuentran dirigidos de una variable a otra, indicando, de esta forma, la existencia de dependencia entre ambas. Los modelos más comunes, y objeto de este trabajo, son los grafos que carecen de caminos cerrados en su estructura, los grafos dirigidos acíclicos, *DAGs* de sus siglas en inglés (*Directed Acyclic graphs*) o Redes Bayesianas.
2. **Modelos gráficos no dirigidos** o también conocidos como *Markov Random Fields*. Como se verá más adelante, las Redes Bayesianas corresponden a un tipo especial de factorización de una distribución de probabilidad conjunta, en la cual cada factor es una distribución por si mismo, mientras que *Markov Random Fields* representan una factorización alternativa, basada en factores o potenciales.

### 2.3.2. Redes Bayesianas

Como se ha visto antes, las Redes Bayesianas son modelos gráficos probabilísticos dirigidos, en concreto, DAGs.

En un grafo dirigido, definimos una relación padre/hijo cuando existe un enlace entre dos variables. De esta forma, en la Red Bayesiana presentada en la figura 2.1, podemos decir que  $X_1$  es el nodo padre de  $X_2$ , que  $X_2$  es padre de  $X_3$ , y, a vez que  $X_3$  es hijo de  $X_2$  y  $X_2$  hijo de  $X_1$ .

Si aplicamos la regla de la cadena (2.5) para un problema genérico de 3 variables  $X_1$ ,  $X_2$  y  $X_3$ , tenemos que:

$$\begin{aligned} P(X_3, X_2, X_1) &= P(X_3|X_2, X_1)P(X_2, X_1) \\ &= P(X_3|X_2, X_1)P(X_2|X_1)P(X_1) \end{aligned} \tag{2.10}$$

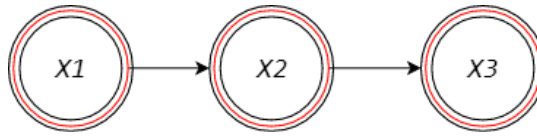


Figura 2.1: Ejemplo de Red Bayesiana

Ahora bien, asumiendo que la Red anterior define el conjunto de independencias condicionales del problema, obtenemos directamente del grafo una simplificación de la expresión anterior:

$$(X_3 \perp\!\!\!\perp X_1) | X_2 \rightarrow P(X_3, X_2, X_1) = P(X_3 | X_2) P(X_2 | X_1) P(X_1) \quad (2.11)$$

de esta forma, una Red Bayesiana puede ser parametrizada a partir de todas las distribuciones de probabilidad condicionales o CPD (*Conditional Probability Distributions*)  $P(X_i | Pa_i)$ , donde  $X_i$  representa el nodo  $i$  y  $Pa_i$  sus nodos padre. De esta forma, generalizando la ecuación 2.9, se define la relación entre un grafo y la función de probabilidad conjunta de un problema mediante la *Regla de la Cadena* del grafo:

La distribución conjunta definida por un grafo es dada por el producto, para todas las variables del grafo, de las probabilidades condicionales de cada nodo con respecto a sus padres (CPDs). Teniendo, para un grafo de  $K$  nodos:

$$p(x) = \prod_{k=1}^K p(x_k | pa_k) \quad (2.12)$$

Esta expresión define la propiedad de *factorización* de la probabilidad conjunta de un problema dado el conjunto de independencias condicionales definidas por una Red Bayesiana.

Una vez introducidas las Redes Bayesianas, en las siguientes secciones se realiza una introducción a la inferencia y aprendizaje en éstas, centrándose en Redes Bayesianas Gaussianas, Redes en las cuales los nodos representan variables aleatorias continuas con distribución Gaussiana.

### 2.3.3. Inferencia en Redes Bayesianas

Se denomina *Inferencia* al proceso de cálculo de la probabilidad una vez conocidos los valores que toman otras variables de la red, es decir, cuando se introduce una determinada *evidencia*.

Así pues, sea  $O$  el conjunto total de variables en la red,  $Q$  el conjunto de *variables observadas*, de las cuales se introduce una evidencia ( $Q = \mathbf{q}$ ),  $R$  el conjunto de variables *no observadas o latentes* dentro del cual, se encuentra el subconjunto  $S$  que incluye las variables a inferir y  $T$  las variables carentes de interés. Cumpliéndose:  $O = Q \cup R$  y  $R = S \cup T$ .



Podemos definir el proceso de inferencia como el cálculo de:

$$\begin{aligned}
 P(S|Q = \mathbf{q}) &= \frac{P(S, Q = \mathbf{q})}{P(Q)} \\
 &= \frac{\sum_T P(S, Q = \mathbf{q}, T = t_i)}{P(Q)} \\
 &= \frac{\sum_T P(S, Q = \mathbf{q}, T = t_i)}{\sum_R P(Q, R = r_j)}
 \end{aligned} \tag{2.13}$$

Así pues, como se puede ver en 2.11, el proceso de inferencia requiere el cálculo de dos marginalizaciones, lo cual, para modelos complejos, supone un alto coste computacional. Para la realización de la inferencia de forma eficiente existen numerosos algoritmos que se pueden dividir en:

1. **Inferencia exacta:** Los casos en los cuales se puede hacer una inferencia exacta son limitados, en especial, redes con todos los nodos latentes discretos o con distribuciones gaussianas univariadas. En este último caso la inferencia se puede realizar algebraicamente al ser la red una parametrización de una distribución *Gaussiana Multivariada Conjunta*. Los algoritmos utilizados en este caso se basan en “empujar las sumas” en el cálculo de la inferencia, es el caso del *Variable Elimination* [3], o en el paso de mensajes a través de árboles, como el algoritmo *Junction Tree* [4] [5].
2. **Inferencia aproximada:** Incluso en los casos en los que la inferencia exacta es posible, puede ser computacionalmente demasiado lenta debido a la complejidad de la estructura de la red o a la complejidad de las distribuciones continuas de los nodos. Se utilizan técnicas de inferencia por *Monte-Carlo*, *Loopy Belief Propagation*, *inferencia variacional*, etc. En este documento no se realizará un estudio de este tipo de técnicas.

## Inferencia en Redes Bayesianas Gaussianas

Tal y como se ve en [6], debido a las características estadísticas de la función normal la inferencia en este caso se puede calcular matricialmente de una forma analítica:

Sea  $X$  el conjunto de  $N$  variables aleatorias del modelo, de las cuales  $q$  son latentes, gaussianas multivariadas a inferir. Podemos particionar  $\mu$  y  $\Sigma$  tal que:

$$\begin{aligned}
 X &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \text{de tamaño } \begin{bmatrix} q \times 1 \\ (N - q) \times 1 \end{bmatrix} \\
 \mu &= \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \text{de tamaño } \begin{bmatrix} q \times 1 \\ (N - q) \times 1 \end{bmatrix} \\
 \Sigma &= \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \text{de tamaño } \begin{bmatrix} q \times q & q \times (N - q) \\ (N - q) \times q & (N - q) \times (N - q) \end{bmatrix}
 \end{aligned} \tag{2.14}$$

Entonces, dada la evidencia  $a$ ,  $P(x_1|x_2 = a)$  es una multivariada definida por:

$$\begin{aligned}
 \bar{\mu} &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(a - \mu_2) \\
 \bar{\Sigma} &= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}
 \end{aligned} \tag{2.15}$$

### **Variable Elimination**

Como define [4], debido a la estructura de la red bayesiana, varias subexpresiones dentro de la expresión de la probabilidad conjunta (ver expresión 2.10) dependen de pocas variables. Es por esto que, calculando dichas expresiones una vez y guardando los resultados, se puede evitar recalcularlas un número exponencial de veces.

Así pues, sea la Red Bayesiana descrita por:

$$P(A, B, C, D) = P(A)P(B|A)P(C|B)P(D|C) \quad (2.16)$$

La probabilidad marginal  $P(D)$  se calcula, en el caso en el que las distribuciones de los nodos es continua (nótese que si fuese discreta se reemplazarían las integrales por sumatorios):

$$P(D) = \int_A \int_B \int_C P(A)P(B|A)P(C|B)P(D|C) \quad (2.17)$$

Empujando las integrales, clave del algoritmo *Variable Elimination*, obtenemos:

$$\begin{aligned} P(D) &= \int_C P(D|C) \int_B P(C|B) \int_A P(A)P(B|A) \\ &= \int_C P(D|C) \int_B P(C|B)\tau(B) \\ &= \int_C P(D|C)\tau(C) \end{aligned} \quad (2.18)$$

Conllevando de esta forma una reducción en el coste computacional en el cálculo de la probabilidad marginal, clave para la inferencia. El algoritmo *Variable Elimination* se puede considerar como un algoritmo de paso de mensaje, en el cual, tras eliminar una variable marginalizándola, se envía el mensaje  $\tau(\cdot)$  al siguiente nodo a eliminar.

#### **2.3.4. Aprendizaje en Redes Bayesianas**

En una red Bayesiana distinguimos dos tipos de aprendizaje: Aprendizaje de parámetros o *model fitting* o aprendizaje de la estructura también llamado *model selection*. Al ser fijada la estructura del modelo con la ayuda del conocimiento de *La Empresa*, este trabajo se centrará únicamente en el aprendizaje de los parámetros del modelo, es decir, las *CPDs* de cada uno de los nodos de la red.

Así pues, el aprendizaje de la red consiste únicamente en ajustar los parámetros que definen la función de probabilidad de cada nodo, de forma que se describa el comportamiento estadístico de los datos. Dentro del aprendizaje de parámetros de la red, encontramos dos casos diferentes.

1. **Caso completamente observado:** Caso en el que no hay datos perdidos en el conjunto de entrenamiento. Si el modelo es una Red Bayesiana Gaussiana, las CPDs, tal y como se determina en [7], pueden ser estimadas independientemente de cada una mediante *Maximum likelihood estimation* (MLE).
2. **Caso parcialmente observado:** Caso en el que hay datos perdidos en el conjunto de entrenamiento. Para el aprendizaje en este caso, el algoritmo más común es el *Expectation-Maximization* (EM) [15].

## Técnica de Máxima Verosimilitud

Como veremos más adelante, en este trabajo nos centraremos en el uso de la técnica de máxima verosimilitud, descrito en [7][8], para el entrenamiento de los parámetros de la red.

Así pues, sea  $X = \{x_1, \dots, x_M\}$  un conjunto de observaciones independientes de una misma variable, con una función de distribución, de la que únicamente se sabe la familia a la que pertenece, siendo  $\theta$  el conjunto de parámetros que la definen. El algoritmo MLE presende asignar un valor a los parámetros  $\theta_X$  que maximice la verosimilitud de los datos con el modelo. Es decir, busca un modelo que sea consistente con los datos, desde el punto de vista de su verosimilitud.

Para ello, se siguen los siguientes pasos:

1. Se parte de la función de densidad de probabilidad fijando los valores observados como fijos. Esta función es llamada *función de verosimilitud*:

$$L(\theta, x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i|\theta) \quad (2.19)$$

2. Se calcula el logaritmo:

$$\hat{l}(\theta|x_1, x_2, \dots, x_n) = \sum_{i=1}^n \log(f(x_i|\theta)) \quad (2.20)$$

3. Se maximiza la función 2.18

$$\hat{\theta}_X = \arg \max \hat{l}(\theta|x_1, x_2, \dots, x_M) \quad (2.21)$$

De esta forma, se obtienen los parámetros  $\hat{\theta}_X$  más cercanos a los parámetros  $\theta$  que definen la distribución real de los datos. Dichos parámetros, pueden ser obtenidos analíticamente a partir de los datos o, por lo contrario, necesitar aproximaciones.

## Máxima Verosimilitud en distribuciones gaussianas multivariadas

Como ya ha sido mencionado anteriormente, en el caso en el que todos los nodos de la red correspondan a variables con distribución gaussiana univariada, la red es una simple representación de una distribución Gaussiana Multivariada, cuyos parámetros que la definen son un vector de medias,  $\mu$ , y una matriz de covarianzas  $\Sigma$ , tal que:

$$f(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp -\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu) \quad (2.22)$$

Siguiendo el desarrollo previamente definido, tal como se describe en [8], obtenemos los valores de la matriz de covarianza  $\Sigma$  y  $\mu$  son:

$$\begin{aligned} \hat{\mu} &= \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{\Sigma} &= \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T \end{aligned} \quad (2.23)$$

Como se puede observar en las ecuaciones definidas en 2.21, los valores de los parámetros que definen la distribución Gaussiana Multivariada se pueden extraer directamente a partir de los datos de entrenamiento. Tras ello, se podrá realizar la inferencia exacta utilizando las técnicas analíticas antes descritas.

## 2.4. Introducción a las Redes Bayesianas Dinámicas

Un gran conjunto de eventos no pueden ser descritos únicamente a partir de un instante temporal. Dado que una Red Bayesiana no representa ninguna dependencia del modelo con respecto al tiempo, surgen las denominadas Redes Bayesianas Dinámicas.

De esta forma, las *Redes Bayesianas Dinámicas* [9][10][11] o, de sus siglas en inglés, *DBNs* (*Dynamic Bayesian Networks*) son extensiones de las Redes Bayesianas vistas anteriormente, utilizadas para la descripción de procesos dinámicos, como las señales que se consideran en este TFM.

Pese a que el concepto de modelo dinámico puede ir intuitivamente asociado a un modelo que varía a lo largo del tiempo, este no es el caso de las Redes Bayesianas Dinámicas, en las cuales el concepto de dinamismo va directamente asociado a la existencia de dependencias con respecto a instantes temporales anteriores al actual.

Así pues una DBN consiste en un conjunto de particiones que representan el estado de las variables en un determinado instante de tiempo. Como podemos ver en la figura 2.2, para cada instante de tiempo, se define una estructura de dependencias entre las variables en ese determinado tiempo. Esta red se denomina la *Red Base* de la Red Bayesiana Dinámica, que, en el caso de la figura ejemplo, está definida por los enlaces azules.

Además de las dependencias dentro del mismo instante temporal, hay enlaces entre variables de diferentes instantes temporales, al conjunto de estas dependencias se las denomina la *Red de Transición*, que se encuentra definida mediante enlaces rojos en la figura.

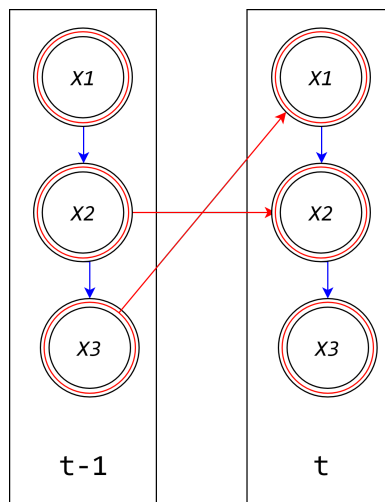


Figura 2.2: Ejemplo de Red Bayesiana Dinámica.

Nótese que la figura 2.2 representa una red en la cual, las variables en un instante  $t$  dependen únicamente de las variables en un instante  $t$  y  $t - 1$ . Sin embargo, es posible que los valores de las variables dependan de varios instantes temporales previos, añadiendo una gran complejidad a la red y, por lo tanto, siendo necesarios algoritmos de aprendizaje e inferencia óptimos computacionalmente.

En este trabajo no se profundizará en estas técnicas debido a la utilización de una Red Bayesiana Gaussiana poco compleja como base de la DBN propuesta y la utilización de una Red de Transición simple, en la que las variables en un instante  $t$  solo dependerán del estado de las variables en un único instante temporal previo. Es por esto que este trabajo no pretende más que una aproximación inicial al dinamismo en Redes Bayesianas.

## 2.5. Gaussianización de datos

Como se verá más adelante, los datos disponibles no tienen una distribución gaussiana. Por lo tanto, para una mejor descripción de los datos a partir del modelo gráfico probabilístico gaussiano es necesaria una modificación de su distribución, de tal forma que estos pasen a tener una distribución gaussiana. A este proceso se le denomina *Gaussianización* de datos.

A lo largo de este proyecto se ha planteado la utilización de dos técnicas de gaussianización simples para variables unidimensionales: Gaussianización basada en ecualización de histogramas y *Feature Warping*.

### 2.5.1. Gaussianización basada en ecualización de histogramas

Tal y como se describe en [12] [13], sea un vector  $X$  de variables aleatorias con una función de distribución de probabilidad conjunta  $f(X)$ . Suponemos que cada variable aleatoria  $x_i$  tiene una función de distribución  $f(x_i)$  y una correspondiente función de distribución acumulada o CDF (*Cummulative Distribution Function*)  $F(x_i)$ . Se supone además que  $\phi(\cdot)$  es la cdf de una variable Gaussiana unidimensional de media cero y varianza unidad, tal que:

$$\phi(\epsilon) = \int_{-\infty}^{\epsilon} \frac{1}{\sqrt{2\pi}} \exp -\frac{\alpha^2}{2} d\alpha \quad (2.24)$$

se puede demostrar que  $Y = \phi^{-1}(F(x_i))$  es una variable aleatoria de media cero y varianza unidad. Siendo  $\phi^{-1}$  la función inversa de la normal, también denominada *probit*. Mediante la gaussianización de las variables con media y varianza unidad, obtenemos la gaussianización total del conjunto de variables.

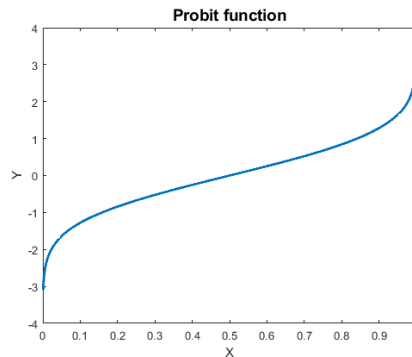


Figura 2.3: Función Probit, o inversa de la normal, para gaussianización de datos con distribución uniforme

### 2.5.2. Feature Warping

La técnica *feature warping* es una técnica utilizada para la modificación de la distribución estadística de una variable a corto plazo, con el fin de que la distribución de la variable completa tenga dicha distribución, en este caso, Gaussiana o Normal.

Para este proyecto, se plantea la implementación propuesta en el artículo [14] para la gaussianización de coeficientes cepstrales. Esta implementación se basa en el siguiente proceso, descrito además en la figura 2.4.

1. Recorrer la señal objetivo mediante una ventana deslizante.
2. Obtener la distribución de los valores de la ventana.
3. Mapeo de la distribución original de la ventana a la distribución deseada, en este caso, distribución Normal de media 0 y varianza unidad, de la siguiente forma:

$$\int_{y=-\infty}^q f(y)dy = \int_{z=-\infty}^m h(z)dz \quad (2.25)$$

donde  $f(\cdot)$  es la distribución original de la ventana deslizante,  $h(\cdot)$  es la distribución gaussiana,  $q$  es el valor original dentro de la ventana deslizante y  $m$  el valor gaussianizado mediante feature warping.

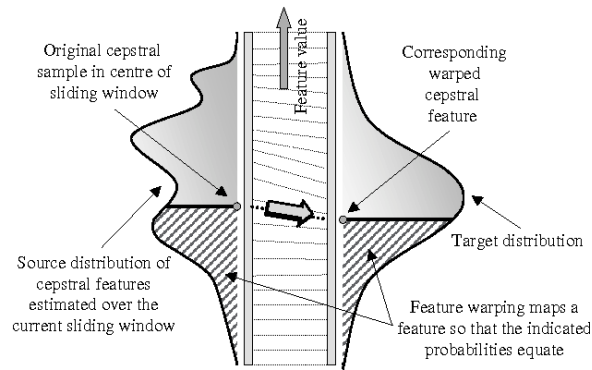


Figura 2.4: Proceso de Gaussianización mediante Feature Warping propuesto en [13]. Esta imagen ha sido extraída del mismo documento

# 3

## Diseño

Tras el estudio de las bases teóricas de este trabajo, en esta sección se describen todas las consideraciones tenidas en cuenta en la etapa de diseño del proyecto.

Esta etapa incluye una descripción de la Base de Datos, que contiene todas las señales de la química del *primario*, y que ha sido proporcionada por *La Empresa*. Además, se realiza una descripción de los sistemas finales propuestos, del método de evaluación utilizado, y del software utilizado para su implementación.

Por último, se realiza una descripción del diseño de la interfaz gráfica de usuario desarrollada para ser utilizada por *La Empresa*, y que le permite, de forma sencilla, la ejecución de los algoritmos desarrollados.

### 3.1. Base de datos

---

La base de datos proporcionada por *La Empresa* se basa en un libro Excel, de 118 669 entradas, con el siguiente formato:

<b>Planta</b>	<b>Ciclo</b>	<b>Fecha</b>	<b>Campo</b>	<b>Medida</b>
<i>Planta 1</i>	22	21/07/2011	<i>Control 1</i>	6.87379
<i>Planta 1</i>	22	22/07/2011	<i>Control 1</i>	6.89552
<i>Planta 1</i>	22	21/07/2011	<i>Control 2</i>	25.08

Donde:

**Planta** es el nombre de una central nuclear gestionada por *La Empresa*. Han sido consideradas 5 centrales nucleares, las cuales denominaremos: *Planta 1*, *Planta 2*, *Planta 3*, *Planta 4* y *Planta 5*. Dichas centrales pueden ser divididas en dos grupos en función de los métodos y protocolos que usan, teniendo por un lado las centrales *Planta 1* y *Planta 2*, y por otro las centrales *Planta 3*, *Planta 4* y *Planta 5*.

**Ciclo** es la numeración que indica el ciclo de trabajo del reactor en el que han sido tomadas las medidas. Como ya fue comentado previamente, denominamos ciclo a la división temporal en

la cual se engloba el funcionamiento del reactor, incluyendo la limpieza y la carga del mismo. El número de ciclos proporcionados por cada central es diferente.

**Fecha** corresponde a la fecha en la que se ha realizado la medida. La periodicidad a la hora de captar una medida depende de la variable a medir y de la central. Un análisis detallado de la frecuencia de captación de muestras será descrito en la sección 4.1.2.

**Campo** es el nombre de la variable cuyo valor ha sido guardado. Pese a que en la base de datos entregada hay un total de 27 campos diferentes, solo serán de interés en este proyecto 10. Dichas variables se pueden dividir en:

- **Señales *Control***: Corresponden a las variables controlables directamente desde cada central. Serán denominados como *Control 1*, *Control 2*, *Control 3* y *Control 4*.
- **Señales *Medida***: Subproductos directamente relacionados con el residuo final en el proceso. Estas señales se encuentran en la base de datos divididas en sus dos diferentes componentes, dependiendo de su solubilidad química (soluble, s; e insoluble, i). Dichas componentes serán denominadas *Medida 1s*, *Medida 1i*, *Medida 2s*, *Medida 2i*, *Medida 3s*, *Medida 3i*, *Medida 4s* y *Medida 4i*.

Por último, el campo Medida es *el valor de la señal “Campo”, adquirido en el ciclo “Ciclo” de la central “Planta”, el día “Fecha”*. Es importante tener en cuenta que una variable puede tener diferentes rangos dinámicos en función de la central, y, que puede haber medidas erróneas debido algún problema en el sistema de captación de las centrales. Un análisis detallado de las señales será descrito en 4.1.1.

## 3.2. Diseño del sistema

---

Para la predicción de las señales *Medida* a partir de las señales *Control*, se proponen dos modelos gráficos probabilísticos para cada una de las centrales. En concreto, se propone una Red Bayesiana Gaussiana y, con el fin de mejorar su rendimiento añadiendo información relativa a las señales de los instantes temporales previos, una Red Bayesiana Gaussiana Dinámica.

Ambos modelos han sido implementados completamente en Matlab<sup>TM</sup> 2018 con el uso del toolbox externo Bayesian Network Toolbox (BNT). BNT [15] es un paquete de Matlab *open-source*, desarrollado entre 1997-2002 por Kevin Murphy, para modelos gráficos dirigidos. BNT soporta una gran cantidad de diferentes distribuciones de probabilidad para sus nodos, diferentes algoritmos de inferencia exacta y aproximada, aprendizaje de parámetros y de la estructura del grafo, además de permitir la implementación de modelos estáticos y dinámicos.

Además, BNT proporciona una amplia documentación con ejemplos incluidos de todo el código disponible, lo que lo convierte en una herramienta muy útil para investigación y para una rápida implementación de prototipos.

Así pues, se proponen los modelos presentados en las figuras 3.1 y 3.2, siendo ambos dos Redes Bayesianas pero, en el primer caso una Red Bayesiana Estática y en el segundo caso una Red Bayesiana Dinámica.



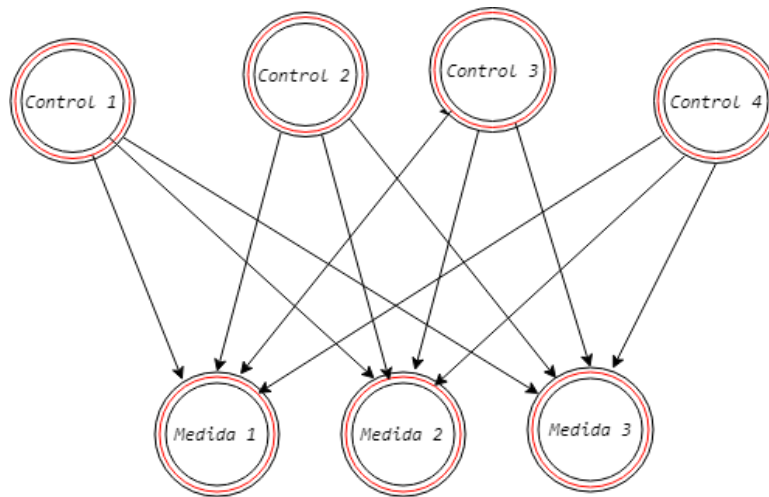


Figura 3.1: Red Bayesiana Final

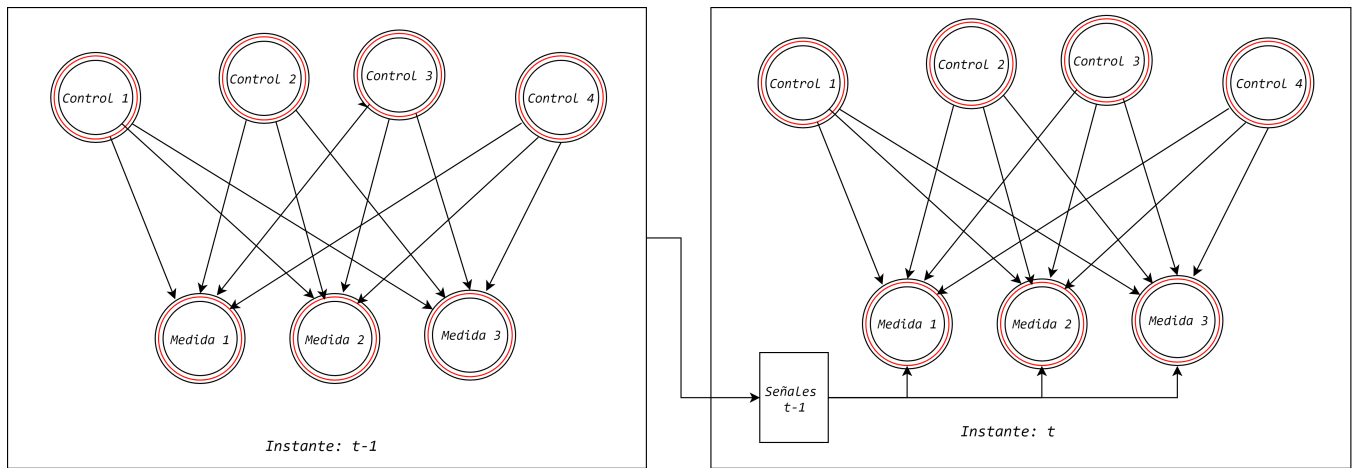


Figura 3.2: Red Bayesiana Dinámica Final

Analizando la estructura de ambos modelos, supondremos, por indicación de *La Empresa* que las variables controlables a lo largo del ciclo, *Control 1*, *Control 2*, *Control 3* y *Control 4*, son independientes entre sí cuando no se conocen las variables medida. Por otro lado, las variables *Medida*, directamente relacionadas con la cantidad de residuo restante en el reactor, y cuyo valor hay que predecir, son: *Medida 1*, *Medida 2*, y *Medida 3*. Estas tres señales se obtienen sumando sus respectivas componentes solubles e insolubles (ver sección anterior) y son independientes entre sí cuando se conocen todas las variables de control. Adicionalmente, todas las variables medida dependen de todas las variables de control.

Adicionalmente, en el caso de la Red Bayesiana Dinámica, vemos que las variables medida en un instante  $t$ , no solo dependen de las variables de control en  $t$ , sino que además dependen del valor de todas las variables en el instante  $t - 1$ . Esta dependencia viene representada en la figura 3.2 por una flecha que relaciona ambos instantes temporales.

En ambos casos, el modelo define siempre una función densidad de probabilidad gaussiana multivariada, y por lo tanto las funciones marginales y las CPDs serán gaussianas, por lo que, previamente al entrenamiento de los modelos mediante MLE (Ver 2.3.4), se realiza una etapa de gaussianización de los datos, tanto en los datos de entrenamiento como los de test, basada en ecualización de histogramas (Ver 2.5.1).

En cuanto a la inferencia de los modelos, ésta es realizada mediante el algoritmo *Variable Elimination*, descrito en la sección 2.3.3 de este documento. Por último, una etapa de *desgaussianización*, inversa a la etapa invertible de gaussianización, es necesaria para la visualización de los datos con sentido químico.

El proceso de diseño se ha realizado en colaboración con *La Empresa*, tal y como se verá en el capítulo 4 de este documento.

### 3.3. Diseño de la evaluación

---

Para la evaluación de los sistemas implementados se propone la utilización de la técnica *K-Fold cross validation* para cada central, siendo  $K$  es el número de ciclos disponibles para dicha central en la base de datos. Así pues, el proceso llevado a cabo es el siguiente para cada una de las centrales:

1. Obtención de los ciclos con valores de todos los elementos del diseño.
2. De los ciclos obtenidos, se utiliza un ciclo para la evaluación del modelo y el resto para su entrenamiento.
3. Se repite 2. de tal forma que se haga predicción de todos los ciclos.

La evaluación de la predicción es realizada mediante dos formas:

1. Visualmente, comparando la predicción de las señales su valor real a lo largo del ciclo de evaluación.
2. Objetivamente, mediante el calculo del *RMSE instantáneo* y *RMSE medio* entre la señal predicha y la señal original. Definiendo el RMSE (*root-mean-square error*) como:

$$RMSE(\hat{\mu}_i) = \sqrt{(\hat{\mu}_i - \mu_i)^2} \quad (3.1)$$

y el RMSE medio como:

$$RMSE = \sqrt{\frac{\sum_{t=1}^T (\hat{\mu}_t - \mu_t)^2}{T}} \quad (3.2)$$

Donde  $\hat{\mu}$  es el valor real de la señal a predecir,  $\mu$  es el valor de la señal predicha por la inferencia en la Red Bayesiana y  $T$  es el número de predicciones realizadas, que, en este caso, es el número de instantes temporales dentro de un ciclo en el que se encuentran las señales del modelo.

Adicionalmente, los resultados de todos los ciclos serán agrupados para cada una de las centrales, calculando el RMSE medio obtenido para cada central.

### 3.4. Diseño de la Aplicación

---

Para el desarrollo de la interfaz de usuario se ha decidido utilizar Matlab<sup>TM</sup>. Las razones por las que se ha decidido su implementación en Matlab<sup>TM</sup> son:

- El toolbox de Matlab<sup>TM</sup> GUI permite la implementación de interfaces de usuario de forma sencilla y rápida. Estas interfaces pueden ser ejecutadas por usuarios con y sin Matlab<sup>TM</sup> en sus terminales.

- La utilización de Matlab<sup>TM</sup> para el desarrollo de la aplicación permite la integración de forma casi directa de todo el desarrollo de Redes Bayesianas ya implementadas con BNT.
- Al ser un enfoque inicial, se ha priorizado la utilización de un lenguaje de programación de interfaces de usuario que, tanto el grupo, como *La Empresa*, ya hubiese utilizado. De esta forma, se asegura el cumplimiento de los hitos establecidos al inicio del proyecto.

Para el diseño de la aplicación se ha tomado como referencia el software gratuito para investigación *Hugin Lite*. *Hugin Lite* es una herramienta que permite la generación de Redes Bayesianas de un máximo de 50 variables, con distribuciones tanto discretas como Gaussianas. Adicionalmente, permite el entrenamiento de los parámetros de la red (con un máximo de 500 valores de entrenamiento) y, por último, permite realizar inferencia sobre todas las variables del modelo no observadas. En la siguiente imagen podemos ver una captura del software (versión para Windows), en ella podemos identificar las diferentes funcionalidades de interés para nuestra propia aplicación.

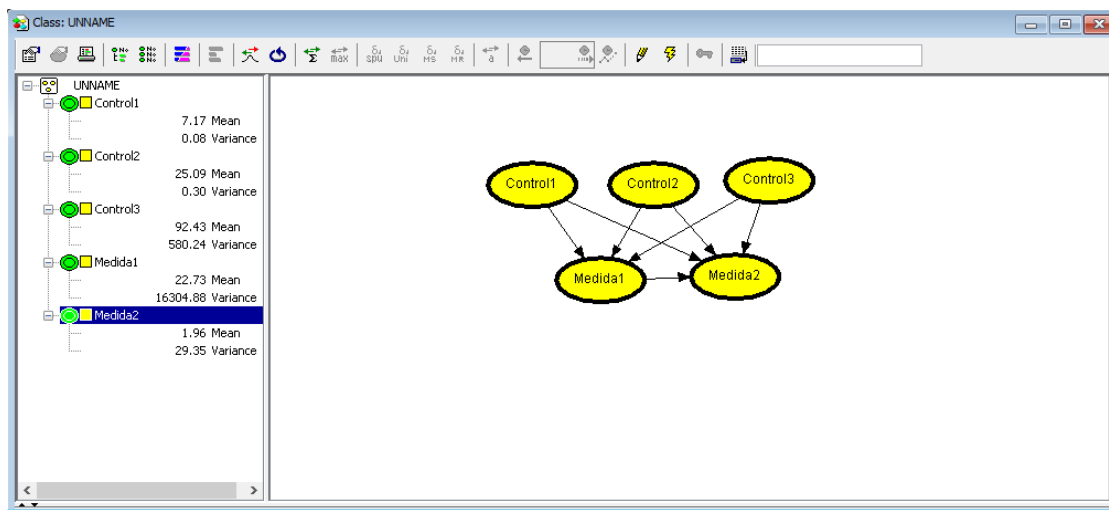


Figura 3.3: Ejemplo Hugin sin evidencias de la Red Bayesiana inicial

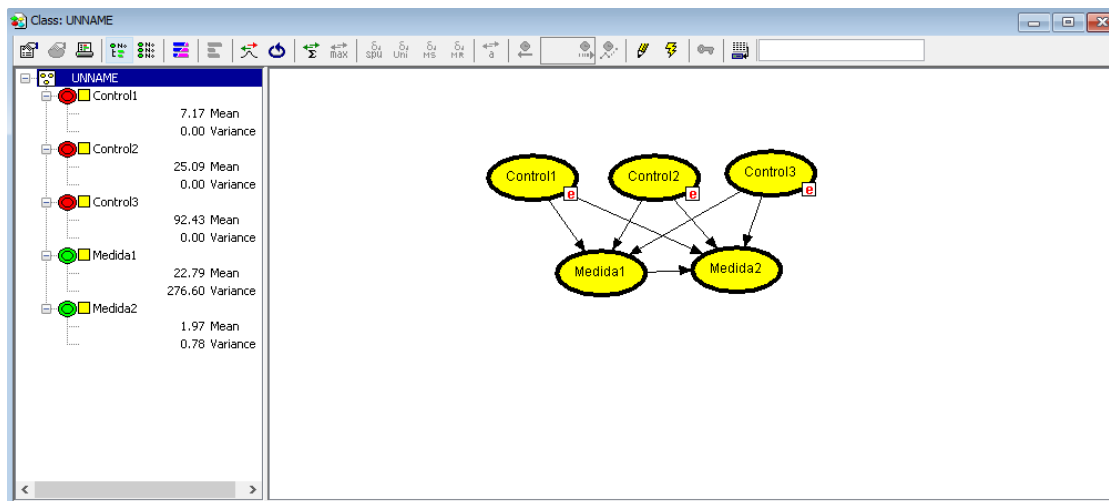


Figura 3.4: Ejemplo Hugin con evidencias de la Red Bayesiana inicial

Podemos ver en ambas figuras que el panel principal de la aplicación se divide a su vez en dos paneles. El de la derecha, permite la visualización y la modificación de la estructura de la

red. En el panel de la izquierda, una vez generada la red, entrenada y compilada, se pueden visualizar los valores predichos por la red, tanto su media como la varianza.

Como se puede ver en la figura 3.4, *Hugin* permite la entrada de evidencias en cualquiera de los nodos, apareciendo estos automáticamente en rojo, a la vez que para el resto de variables, se presentan los valores de media y varianza predichos.

La idea para el desarrollo de la interfaz para *La Empresa* es recrear el funcionamiento de Hugin en Matlab<sup>TM</sup>, salvo por dos diferencias:

1. La estructura del modelo no podrá ser modificada. Lo que supone que *La Empresa* solo podrá utilizar el modelo presentado en este trabajo, sin poder añadir o eliminar ni variables ni dependencias en el modelo.
2. Por interés de la empresa, pese a que la *demo* incluirá un modelo para cada central ya entrenado con los datos proporcionados al principio del proyecto, se pretende que el usuario pueda reentrenar la red a partir de nuevos datos en un fichero Excel con el mismo formato que el ya presentado con anterioridad.

Las especificaciones técnicas del desarrollo de la aplicación serán detalladas en la sección 4.5.

# 4

## Desarrollo del Proyecto

En este apartado se presenta el desarrollo del proyecto, incluyendo, desde el análisis preliminar de las señales de la base de datos, hasta una explicación de todos los pasos realizados para la implementación de los modelos presentados previamente.

### 4.1. Análisis Preliminar

---

El primer paso, una vez iniciado el proyecto con *La Empresa*, es la realización de un análisis de las señales de la base de datos. El análisis realizado puede dividirse en un análisis del comportamiento, de la frecuencia de captación y de la correlación de las diferentes señales de interés.

#### 4.1.1. Análisis del comportamiento de las señales

Como se ha comentado previamente en este documento, encontramos 10 señales de interés para *La Empresa*. De estas, hay 4 señales *Control*, aditivos químicos o parámetros controlables de la central, que se incluyen a lo largo del ciclo de actuación del reactor, y 6 señales que conforman las 3 señales denominadas *Medida*, cuya variación depende, en principio, del valor de las señales de control y que están directamente relacionadas con la liberación de productos de corrosión.

Para el primer análisis de las series temporales se ha tomado como referencia temporal el día 20 de junio de 2009. Esta será nuestra referencia de tiempos, denominada “ $t=1$ ”, siendo los números enteros siguientes la referencia temporal en días (es decir, el 21 de junio de 2009 es “ $t=2$ ”, etc.). El eje de tiempos se considera continuo, para así poder acomodar varias medidas en horas diferentes de un mismo día, calculando el momento exacto de adquisición de la medida con la correspondiente fracción del día de adquisición. Nos referiremos a las unidades en el eje de abscisas de las series como “unidades temporales” o simplemente “ $t$ ” (indicando “2300 unidades temporales” y “ $t=2300$ ” de forma indistinta).

Así pues para analizar el comportamiento de las señales participantes en el modelo, se realizan representaciones de todas las señales, divididas en ciclos, para todas las centrales, tal y como se visualiza en la figura 4.1. Debido a los límites en la extensión de este TFM, se incluyen en esta memoria conclusiones generales y ejemplos de este análisis. Sin embargo, por completitud,

se adjuntan todas las gráficas necesarias para el análisis en el Anexo A. Pese a que la visuación de dicho anexo, al igual que el resto de los anexos, es recomendable, no es necesario para el correcto seguimiento de este trabajo.

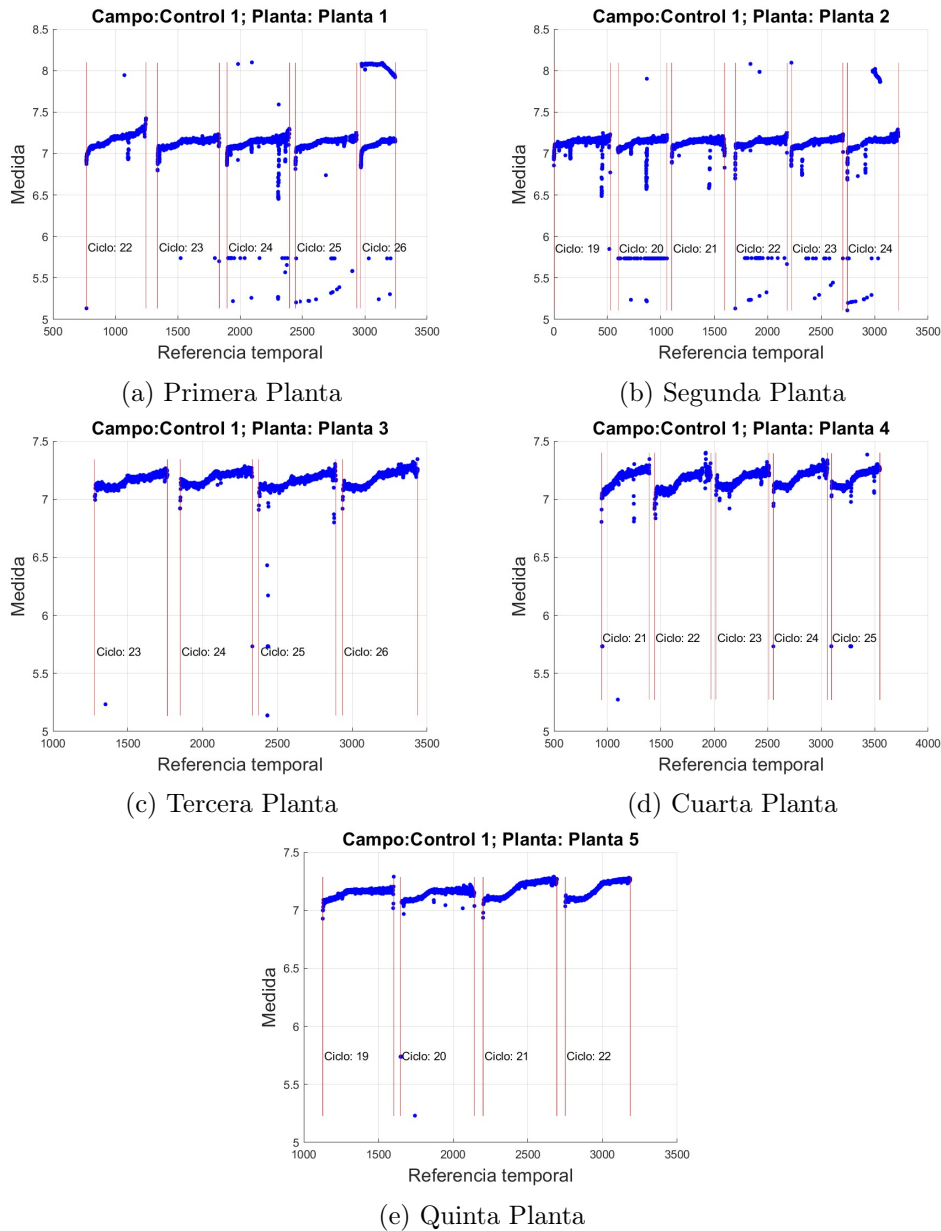


Figura 4.1: Señal *Control 1*

A continuación se expresan las conclusiones obtenidas para cada una de las señales. Adicionalmente, se incluye como ejemplo, en la figura 4.1, la representación temporal de la señal control 1. Para el resto de las señales, se incluyen solamente las siguientes conclusiones:

### Señales de control

#### 1. *Control 1*

En las gráficas de la figura 4.1 podemos ver como la señal *Control 1* presenta valores similares para cada ciclo de cada una de las centrales. La señal, salvo algunos picos, es creciente en el rango comprendido entre 6.5 y 7.5. Fuera de este rango, podemos considerar

que los valores de la señal son atípicos, y *La Empresa* nos indicó que no los tuviéramos en cuenta.

#### 2. *Control 2*

En la figura A.2 vemos como la señal *Control 2*, pese a tener una tendencia casi constante dentro de cada ciclo, posee un valor medio que varía en función de la central que adquiere la medida (véase la diferencia entre los valores de la señal para la planta 1 y para la planta 5), pero también puede variar considerablemente para cada ciclo dentro de una misma central (véase la central 4 y 3).

La constancia de las señales dentro de cada ciclo, y la gran diferencia en la media de la señal para diferentes ciclos y centrales, junto con las indicaciones de *La Empresa* de acuerdo con su conocimiento experto, plantean el uso de la variación de la señal en vez del valor de la misma para el desarrollo de los modelos.

#### 3. *Control 3*

En el caso de la señal *Control 3*, la tendencia es prácticamente constante en todos los ciclos y todas las centrales. Sin embargo, presenta varios picos hacia valores inferiores a la tendencia de la serie, alcanzando valores nulos o casi nulos.

Debido al rango dinámico de la señal, comprendido entre 0 y 100 o entre 0 y 1, se considera que es una señal porcentual cuya tendencia es mantenerse cerca del 100%, salvo por bajadas puntuales de breves instantes temporales.

#### 4. *Control 4*

La señal *Control 4*, es la señal con mayor variación intra-ciclo de todas las señales de control, lo que imposibilita la identificación de una tendencia clara. En todas las centrales el rango dinámico es similar, abarcando la mayor parte de la señal dentro del rango entre 25 y 50.

### Variables Medida

De la misma forma que para las señales de control, en la figura 4.2 se incluye como ejemplo la representación temporal de la señal *Medida 1*, mientras que para el resto de las señales se incluyen solamente las conclusiones extraídas del análisis.

La señal *Medida 1*, al igual que todas las señales *Medida*, es obtenida mediante la suma de sus dos respectivas componentes, *Medida 1s* y *Medida 1i*, cuyos valores son los disponibles en la base de datos proporcionada. Al tener ambos casos una frecuencia de adquisición diferente, y al realizarse la suma de ambas únicamente en instantes temporales en los que ambas medidas están disponibles, el número final de valores disponibles de todas las señales *Medida* es menor que el de las señales de control.

#### 1. *Medida 1*

El comportamiento de *Medida 1* varía considerablemente en función de la central que capta sus valores. En el caso de *Planta 1* y de *Planta 2*, tenemos dos señales altamente influenciadas por valores atípicos de la señal, y que, por consideraciones de *La Empresa*, no se pueden considerar valores erróneos y se tienen que tener en cuenta. En las centrales *Planta 3* y *Planta 4* no se encuentran estos valores atípicos y podemos apreciar como la señal tiende a tener valores máximos en los instantes temporales centrales de cada uno de los ciclos, siendo en ambos casos el rango dinámico bien definido como [0-25].

## 2. Medida 2

El comportamiento de dicha señal es muy similar al de *Medida 1*, en el cual, para las centrales *Planta 1*, *Planta 2* y *Planta 5* la tendencia de la señal se ve altamente influenciada por valores atípicos de una gran magnitud. Sin embargo, el rango dinámico de esta señal es considerablemente inferior al de *Medida 1* para todas las centrales.

Por otra parte, las señales de la *Planta 2* y *Planta 3* están comprendidas en un rango muy limitado, [0-10], con una tendencia muy similar a la encontrada en la señal anterior estando ambas fuertemente correlacionadas.

## 3. Medida 3

Por último, la *Medida 3* tiene un comportamiento similar a las 2 medidas anteriores. Salvo que, en este caso, la tendencia de la señal en la central *Planta 3* también se ve influenciada por valores atípicos altos. En el caso de la central *Planta 4*, la señal se encuentra comprendida dentro del rango [0-18], sin una tendencia clara.

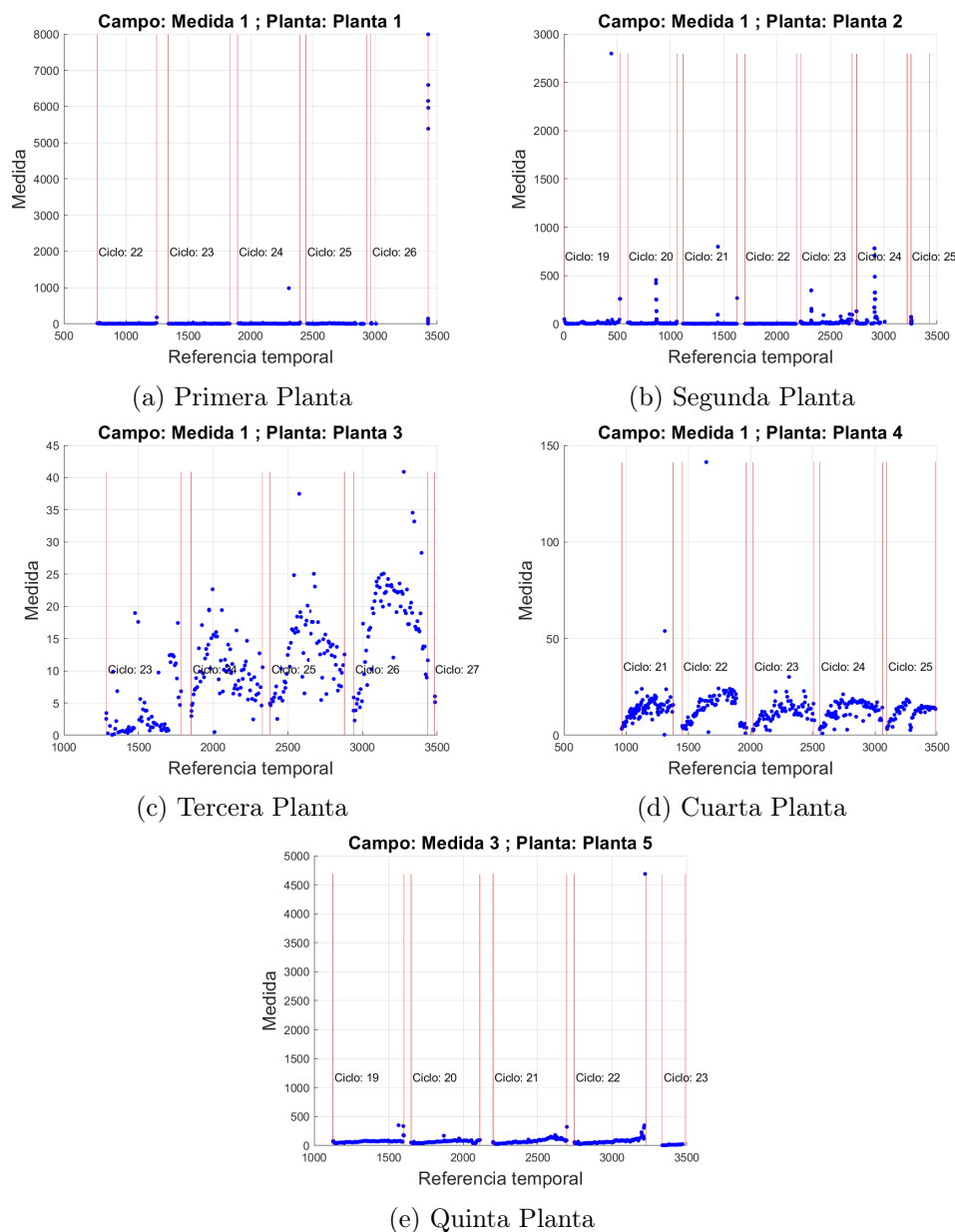


Figura 4.2: Señal *Medida 1*



#### 4.1.2. Análisis comparativo temporal

En este análisis se pretende determinar cuáles son los patrones temporales de adquisición de las señales de cada central y ciclo, con el objetivo de plantear estrategias de homogeneización adecuadas a los datos disponibles.

Cabe destacar que no llamamos al proceso de adquisición de las medidas “muestreo”, pues en tratamiento de señales de este tipo, este tipo de técnicas suelen ir asociadas con patrones de adquisición de medidas uniformes en el tiempo, cosa que no ocurre con las series disponibles en general. Además, las medidas se toman de forma diferente dependiendo de la señal a la que pertenecen. Emplearemos, por lo tanto, la expresión “adquisición”, “medidas adquiridas” o expresiones similares.

Para analizar esto se han realizado las representaciones disponibles en el Anexo ?? del final de este documento. Un ejemplo de éstas representaciones es la figura 4.3. Como se puede observar, en estas gráficas se representan, para cada ciclo disponible de cada central:

- En blanco, instantes temporales con medida para el elemento determinado.
- En negro, instantes temporales sin medida para el elemento determinado.
- Con puntos azules, intervalos temporales entre enteros. Es decir, los puntos azules marcan días con más de una medida.

De esta forma, extraemos como conclusión que, para todas las plantas, la frecuencia de adquisición de las señales *Control* es muy alta, alcanzando para muchos ciclos valores de casi una medida al día. Sin embargo, las señales *Medida* son adquiridas generalmente siguiendo un patrón menos frecuente, con una media de una adquisición cada cinco días.

Además podemos observar que, para las centrales *Planta 1* y *Planta 2* es muy habitual la adquisición de varios valores de señal a lo largo del día, mientras que para el resto de centrales no.

De las gráficas del anexo ?? extraemos además algún caso atípico como:

- La central *Planta 2*, en el ciclo 25, consta de muy pocos valores en la base de datos (tiene una duración de 175 unidades temporales mientras que los demás ciclos tienen una media de 500), además carece de medidas para la señal *Control 1*.
- Para la central *Planta 3*, en el ciclo 27, se tiene demasiadas pocas muestras para ser analizado.
- Para la central *Planta 5*, se tiene que la frecuencia de adquisición de los valores de la variable *Control 4* es muy baja, incluso menor que para la adquisición de variables medida. Esto mismo ocurre para la central *Planta 1*, en el ciclo 26.

Como se verá más adelante, para el desarrollo de los modelos gráficos solo se utilizarán instantes temporales en los que se encuentren valores de todas las señales involucradas, por lo que el número de datos disponibles, tanto para el entrenamiento como para el test de los modelos propuestos, será aproximadamente el número de muestras de las variables *medida* (ya que generalmente los valores de dichas señales se adquieren los mismos días), lo que supone aproximadamente un total de 250 muestras por central.

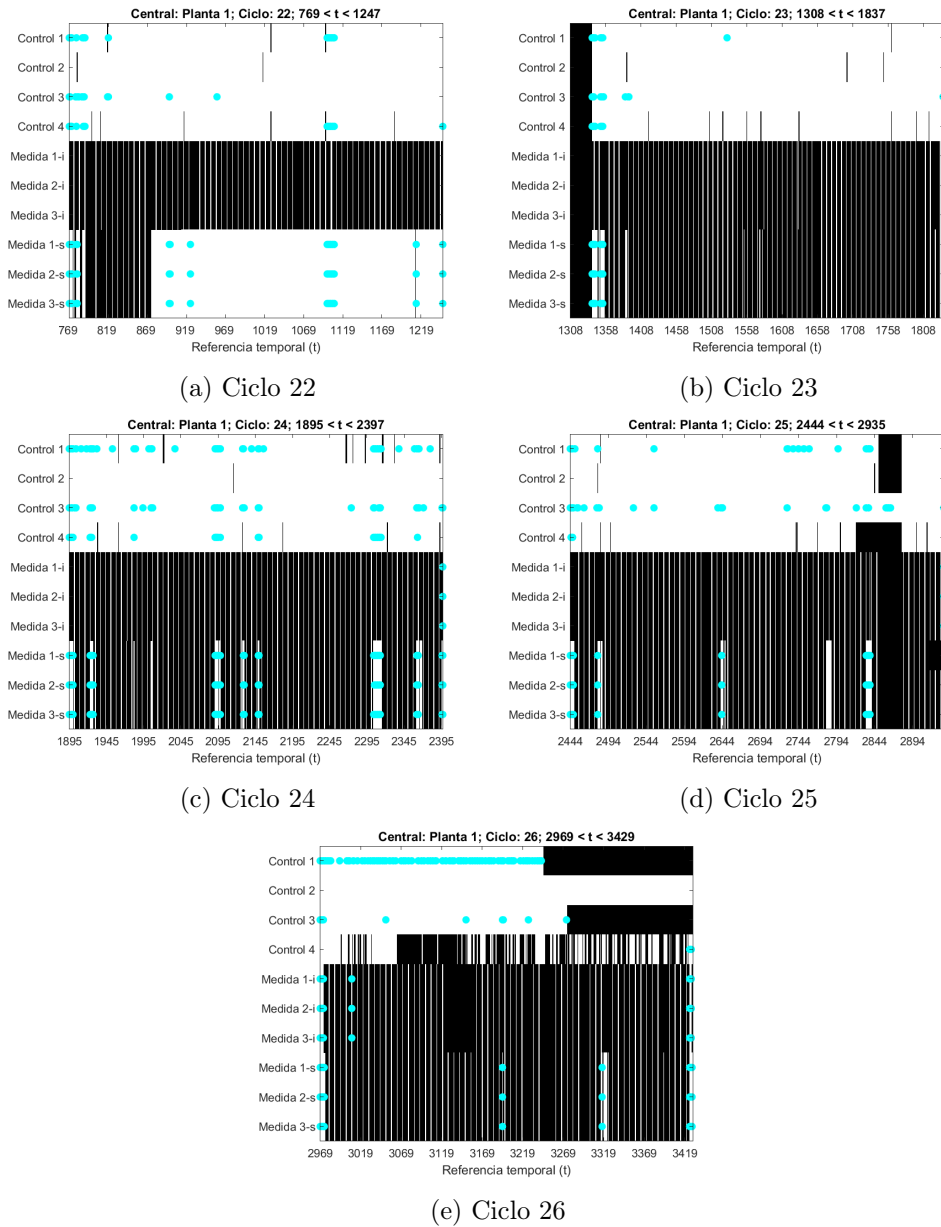


Figura 4.3: Análisis comparativo temporal para la central *Planta 1*

#### 4.1.3. Análisis de la correlación mediante Scatter Plots

Tras los análisis descritos previamente, se decide realizar un análisis de las relaciones de las señales *medida* con las señales *Control*.

La medida de correlación más típica es el coeficiente de correlación de Pearson [16][17]. Este coeficiente puede ser definido como el índice que mide el grado de covariación entre distintas variables relacionadas linealmente. Este índice tiene valores que varían dentro del rango  $[-1:1]$ , cuanto más cerca de 1 y -1 mayor es la correlación, y cuanto más cerca de 0 menor.

Así pues, el coeficiente de correlación de Pearson entre dos variables aleatorias,  $X$  e  $Y$ , se define como:

$$\rho_{X,Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}} \quad (4.1)$$

Donde  $n$  es el conjunto de valores de la variable, y  $\bar{X}$  y  $\bar{Y}$  son las medias muestrales ( $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ ) de los valores de las variables  $X$  e  $Y$ ,

Así pues el análisis propuesto para el estudio de la correlación entre las variables *Control* y las variables *Medida* se basa, tal como podemos ver en la figura 4.4 en:

- Cálculo del coeficiente de correlación de Pearson entre todas las variables *Control* y las variables *Medida* (suma de las dos componentes de cada una de las señales). Esto es el coeficiente de correlación a largo plazo o *long-term*.
- Representación de diagramas de dispersión *Scatter-Plots* que relacionan las variables *Medida* con las variables *control*.
- Correlación a corto plazo. Para ello se recorren las señales a analizar mediante una ventana deslizante de 101 instantes temporales. Se representa dicho coeficiente a lo largo del tiempo.
- Para todo coeficiente de correlación calculado, se realizará un test estadístico mediante una  $t$  de Student, que arrojará un  $p$ -valor. Asumiremos una significancia estadística de 0.05, por lo que  $p < 0,05$  indicará que se puede rechazar la hipótesis nula de que el coeficiente de correlación es nulo; o dicho de otra forma, se rechazará el hecho de que el coeficiente de correlación no es significativo.

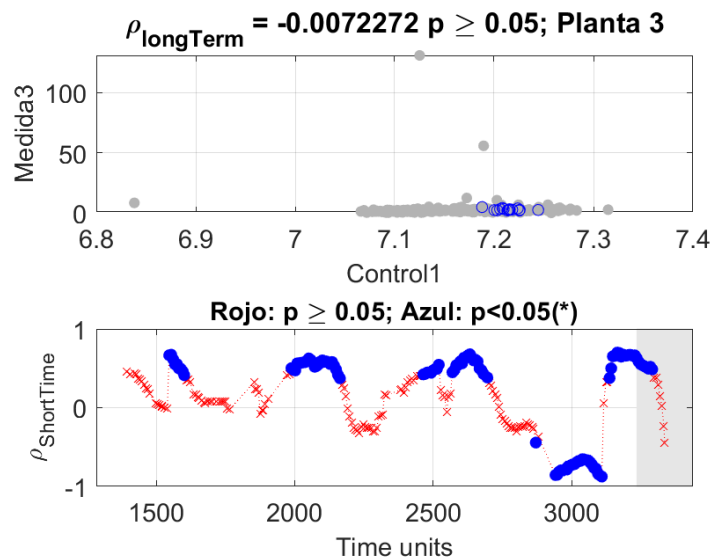


Figura 4.4: Ejemplo del análisis realizado. La gráfica de arriba representa un diagrama de dispersión. La gráfica de abajo representa el coeficiente a corto plazo, en azul si  $p\text{-valor} < 0.05$ , mediante una cruz roja en el caso contrario.

Del análisis realizado, cuyas gráficas están a disposición del lector en el anexo ??, extraemos las siguientes conclusiones:

- Para la central *Planta 1* la correlación entre las variables *Medida* y las *Control* es moderada negativa, al tener de media un valor de  $\rho$  de -0.5. Los máximos de correlación para todas las señales *Medida* se obtiene con la señal *Control 2* alcanzando un valor de  $|\rho| = 0,82$ ,  $|\rho| = 0,81$  y  $|\rho| = 0,75$  para las tres señales de medida respectivamente.
- Pese a que para *Planta 2* la correlación entre las variables disminuye, se puede decir que en general es también moderada negativa, con un valor de  $|\rho|$  medio de 0.42. Sin embargo,

en este caso los máximos de correlación se obtienen con *Control 3* para la *Medida 2* y *Medida 3*, obteniendo valores de  $\rho$  iguales a  $-0,76$  y  $-0,69$  respectivamente, y  $-0,45$  con *Control 4*.

- Por otro lado, en las centrales *Planta 3*, *Planta 4* y *Planta 5* se obtiene generalmente una disminución en la correlación en todas las relaciones entre las variables. Obteniendo en muchos casos valores como los presentes en la figura 4.3. De tal forma que no se encuentra una relación moderadamente lineal entre variables, exceptuando la relación entre las señales *Control 1* y *Medida 1* para la quinta planta.

## 4.2. Primera aproximación

---

Una vez conocidas las señales de interés, se comienza con la generación del modelo gráfico probabilístico que permitirá la predicción de las señales *Medida*, mediante el paquete BNT de Matlab<sup>TM</sup>.

El modelo inicial desarrollado y que, a partir de los sucesivos cambios posteriormente descritos, dará lugar a la red definida en la sección 3.2, es una Red Bayesiana Gaussiana tal como se presenta en la figura 4.5.

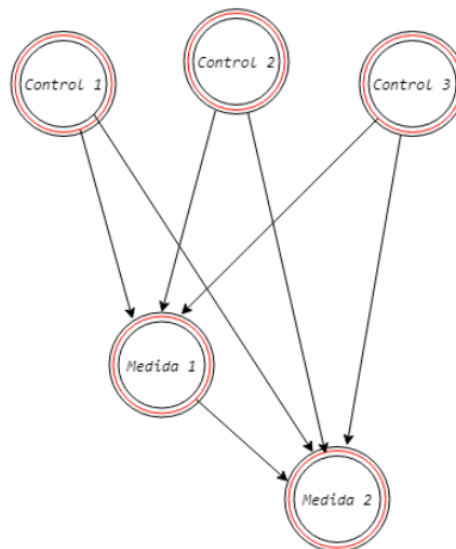


Figura 4.5: Red Bayesiana Inicial

El proceso completo para el desarrollo del modelo queda definido por el diagrama de la figura 4.6, en el cual encontramos los siguientes módulos:

1. **Extracción de los elementos de la base de datos.** De las 27 señales disponibles en la base de datos, se descartan las señales no consideradas de interés. De las señales presentes en el modelo, solo se extraen los instantes temporales en los que han sido adquiridas todas las señales y se realiza la suma de ambas componentes (soluble, s; insoluble, i) de cada señal *Medida*. La razón por la cual solo se capturan instantes temporales en los que se encuentran todas las señales es para evitar errores al realizar *data imputation* de series con valores atípicos, como las presentes en la base de datos. Adicionalmente, para cada

ciclo de cada central, se eliminan los datos provenientes de los primeros y últimos 30 días, al considerarse el comportamiento del reactor inestable, tal y como *La Empresa* nos ha indicado en varias ocasiones.

2. **Entrenamiento de la Red:** Se realiza el entrenamiento de los parámetros que definen el modelo mediante MLE.
3. **Inferencia:** Tras ajustar el modelo a los datos, la inferencia se realiza mediante el algoritmo *Variable Elimination* al ser más óptimo computacionalmente que la inferencia analítica, con la que también se han realizado pruebas.

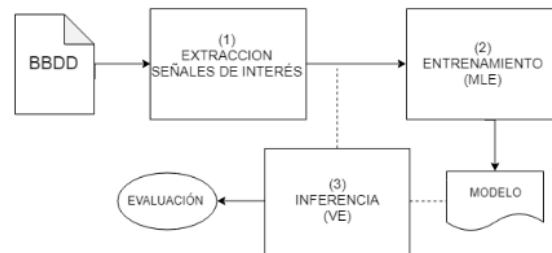


Figura 4.6: Desarrollo del modelo inicial

Los resultados obtenidos para este modelo inicial se describen en la sección 5.1.

### 4.3. Modificación del modelo y de los datos

---

Tras la obtención de los primeros resultados de predicción, y tras el seguimiento del progreso del proyecto por parte de *La Empresa*, se decide realizar una serie de cambios en el modelo y en los datos. En este apartado se describen todos los cambios sucesivos realizados hasta la obtención del modelo final.

#### 4.3.1. Modificaciones directamente sugeridas por *La Empresa*

Durante el proyecto, se realizaron los siguientes cambios a petición de *La Empresa*:

1. Eliminación de la dependencia de la variable *Medida 2* con respecto a *Medida 1*, pues carece de significado químico al provenir ambas señales de reacciones diferentes.
2. Limpieza de los valores atípicos de la variable *Control 1*, tal y como se describió en el análisis preliminar de este capítulo, mediante la eliminación de los valores por encima de 7.5 y por debajo de 6.5, pues estos si que son considerados erróneos.
3. Tal y como ha sido comentado anteriormente, en vez del propio valor de la variable *Control 2* se hará uso de su derivada para así considerar la variación de ésta.
4. Por último, *La Empresa* propone la inclusión en el modelo de una variable de control adicional, *Control 4*, con la intención de predecir mejor las variables *Medida*. Se incluye además la variable *Medida 3* al considerar su utilidad para predicción del residuo final.

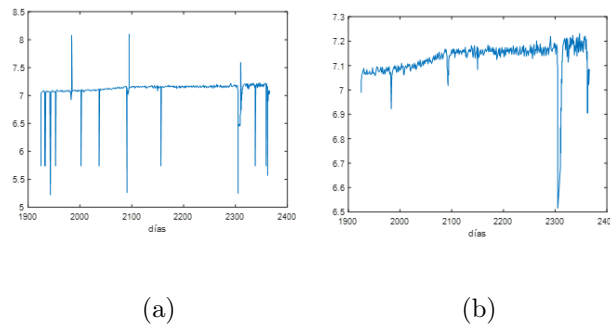


Figura 4.7: Limpieza de *Control 1* para *Planta 1* el ciclo 24. (a) Señal antes del proceso. (b) Señal después del proceso.

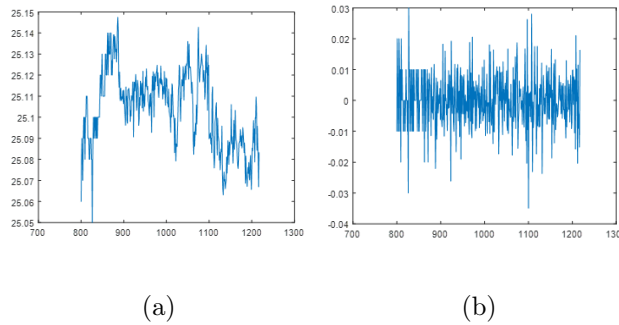


Figura 4.8: Derivación de *Control 2* en *Planta 1* el ciclo 22. (a) Señal inicial. (b) Derivada

De esta forma se obtiene la Red Final ya descrita en la sección 3.2. Adicionalmente, para mejorar el rendimiento de la red, se propone la transformación de los datos de entrada, mediante el uso de las técnicas descritas a continuación.

#### 4.3.2. Gaussianización de los datos

En el diseño de la red se estableció que los diferentes nodos del modelo se definiesen con una distribución Gaussiana, cuyos parámetros, media y varianza eran aprendidos mediante el algoritmo MLE a partir de los datos de entrenamiento disponibles.

Sin embargo, las variables vistas anteriormente no tienen necesariamente una distribución gaussiana. Por lo tanto, se estima que mediante una gaussianización de los datos, los datos transformados se ajustarán mejor al modelo.

La Gaussianización implementada, descrita en la sección 2.5, se basa en ecualización de histogramas de las señales. Los pasos para el proceso se ven reflejados en la figura 4.9:

1. **Aumento de los datos:** Se añaden datos por debajo y por encima de los datos originales para obtener una estabilidad algorítmica y así evitar valores nulos o infinitos.
2. **Ecualización de histograma:** Se aplica sobre los datos la transformación definida por la CDF de la propia variable. El resultado es la uniformización de la distribución de la variable.
3. **Gaussianización:** Por último, se realiza la transformación definida por la función *Probit*, inversa de la función normal, para la gaussianización de los histogramas uniformes.

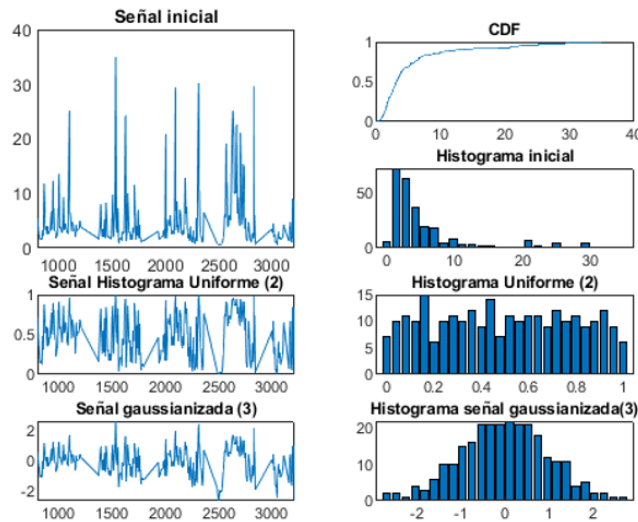


Figura 4.9: Proceso de Gaussianización en señal *Medida 1, Planta 3* ciclo 24

Cabe destacar que, para la correcta inferencia de las variables *Medida* a partir de una evidencia/observación de las variables de control, es necesario, primero una gaussianización de las evidencias de la misma forma que se hizo para las señales en el entrenamiento. Dicha gaussianización será realizada con la CDF obtenida en el entrenamiento.

Además, para la obtención de la predicción final, con significado químico, se requiere una etapa de *desgaussianización*, siguiendo el proceso inverso al descrito anteriormente. La necesidad de esta etapa de desgaussianización es la principal razón por la que se ha realizado gaussianización basada en ecualización de histogramas en vez de *Feature Warping*, al no permitir esta la realización de un proceso inverso.

Tanto para el proceso de desgaussianización, como para la gaussianización de las evidencias, se requiere el almacenamiento de las CDFs de los datos de entrenamiento, para su posterior uso.

#### 4.3.3. Interpolación de datos

Debido a la inclusión de dos parámetros nuevos a la red, el hecho de que solo sean utilizados instantes temporales con muestras de todas las señales, la eliminación de los valores erróneos de *Control 1*, además de la baja frecuencia de captación de las variables medida, se plantea la interpolación de los datos de entrada para el entrenamiento de la red.

Pese a que en un principio se plantea la interpolación de todas las señales de modelo para el entrenamiento de la red, finalmente se decide interpolar únicamente las señales de control para intentar no incluir el entrenamiento de la red demasiados datos erróneos, debido a la inevitable imperfección de la imputación de datos mediante interpolación.

De esta forma, se ha realizado una interpolación lineal de los datos de entrada para el entrenamiento de la red. En el caso en el que los datos faltantes se encuentren al principio o al final de la señal y, por lo tanto, no se pueda realizar la interpolación lineal, se realiza la interpolación al vecino más próximo.

De esta forma, el diagrama que describe el proceso completo desarrollado, antes descrito por 4.10, pasa a ser:

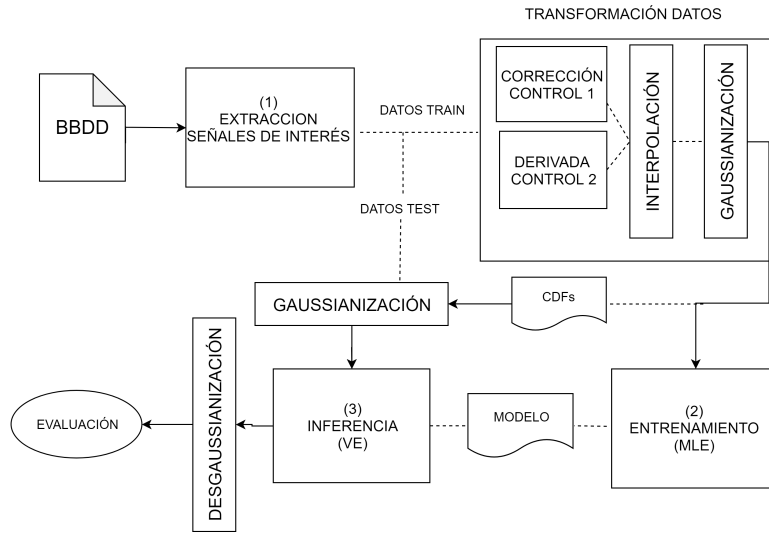


Figura 4.10: Diagrama final del proceso de desarrollo

#### 4.4. Aproximación Dinámica

Tras el desarrollo de la Red Bayesiana previamente descrita. Se realiza el desarrollo de la Red Bayesiana Dinámica presente en la figura 3.2, teniendo en cuenta las siguientes consideraciones:

1. La *Red Base* de la DBN corresponde con la Red Bayesiana previamente implementada. Siguiendo la estructura de dicha red, el valor de las variables latentes en un instante  $t$  dependen únicamente de las variables de control en dicho instante.
2. Adicionalmente, las variables latentes dependen del valor de todas las variables en un instante de tiempo anterior  $t - N$ . Tal como se presenta en la sección 5.3, se ha realizado pruebas con diferentes valores de  $N$ .
3. Aprovechando dicha estructura temporal simple, en la cual solo se tiene en cuenta un instante temporal anterior, se aplican los métodos de inferencia y aprendizaje previamente descritos, que son aptos para grafos de complejidad reducida.
4. El aprendizaje de la red se realiza a partir de los datos Gaussianizados y con interpolación de las variables de control.

#### 4.5. Desarrollo de la Aplicación

La aplicación desarrollada integra la Red Bayesiana Estática presentada anteriormente. Para ello, se ha entrenado un modelo para cada una de las centrales de la base de datos. Cada modelo ha sido entrenado con los datos transformados según la sección anterior, gaussianizados e interpolados de todos los ciclos disponibles para cada central.



De esta forma, la aplicación incluye la gaussianización de los datos de entrada (evidencias), la inferencia mediante el algoritmo *variable elimination*, la desgaussianización de los valores predichos y, por último la visualización de los resultados.

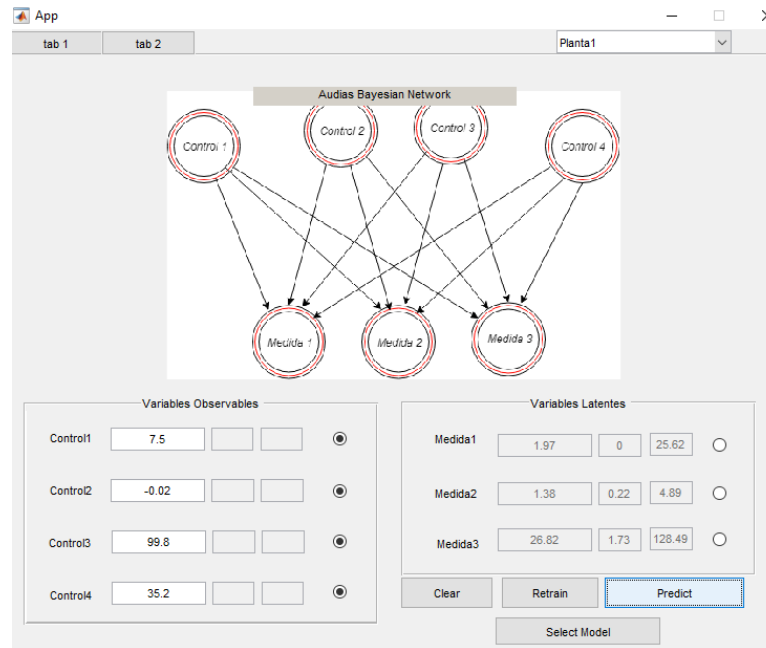


Figura 4.11: Aplicación AUDIAS para el proyecto

Tal y como se ve en la figura 4.11 , la aplicación tiene una estructura inspirada en la de Hugin en la cual encontramos:

1. Panel de visualización de la estructura de la Red: Se recuerda que dicha red no es modificable.
2. Panel de variables: En este panel se deben incluir los valores de las variables marcadas como evidencias / observadas. Y, tras la inferencia (Botón predict), se muestran las predicciones (media y rango de confianza) de las variables no observadas.
3. Botón *Predict*: Realiza la inferencia de todas las variables no observadas.
4. Botón *Clear*: Limpia todos los valores actuales de las variables.
5. Botón *Retrain*: Pese a que la red ya está entrenada con los datos de la base de datos, este botón permite la lectura de un fichero xlsx con formato similar al presentado en 4.1.1, añadiendo una columna *Referencia Temporal* que indica la diferencia de días entre la captura de la medida y la fecha de Referencia.

Planta	Ciclo	Fecha	Ref. Temp	Campo	Medida
Planta 1	22	21/07/2070	769	Control 1	6.87379
Planta 1	22	22/07/2011	770	Control 1	6.89552
Planta 1	22	21/07/2011	769	Control 2	25.08

Y, de esta forma, se entrenará un modelo nuevo a partir de datos nuevos disponibles. El modelo creado, se guardará directamente en formato .mat.

6. Botón *Select Model*: Permite la carga de un modelo ya entrenado.



# 5

## Experimentos y Resultados

En este capítulo se presentan los resultados obtenidos siguiendo el mismo orden que el establecido en el capítulo anterior.

El método de evaluación, descrito en el apartado 3.3, se basa en la utilización, para cada central, de la técnica *K-Fold Cross Validation*, siendo  $K$  el número de ciclos disponibles. De esta forma, la predicción de cada uno de los ciclos, para cada una de las centrales, será evaluada mediante la utilización del RMSE instantáneo (es decir, el RMSE de una predicción en un momento puntual del tiempo) y RMSE medio.

### 5.1. Resultados Red Bayesiana inicial

Para el análisis de las predicciones obtenidas con el primer modelo, se decide la representación de, tal como podemos ver en los ejemplos presentados en las figuras 5.1, 5.2 y 5.3, varias señales divididas en 3 filas diferentes:

- Primera fila: En esta fila se representan los resultados de la predicción a lo largo del ciclo. Dado que el resultado de la inferencia en un instante  $t$  es una función de probabilidad gaussiana, los resultados serán representados a partir del valor medio de dicha distribución,  $\mu_t$ , y su *margen de credibilidad*, que representa la incertidumbre en la predicción, que determinará la incertidumbre en la predicción. Este margen de credibilidad es definido de acuerdo con la mayor parte de la literatura existente ([1]), como:

$$\mu_t \pm 2\sigma_t \tag{5.1}$$

donde  $\mu_t$  y  $\sigma_t$ , son, respectivamente la media y la desviación típica de la predicción. De esta forma, las predicciones son representadas junto con los valores reales de las señales a modo de comparativa.

- Segunda fila: En esta fila se muestran las señales *Control* a lo largo del ciclo.
- Tercera fila: Por último, se representan tanto la medida RMSE en cada instante temporal, como el RMSE medio a lo largo del ciclo predicho.

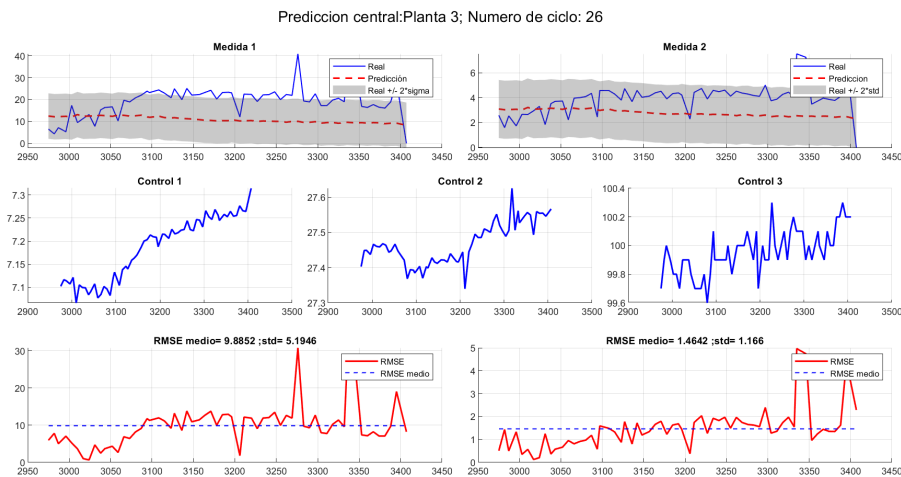


Figura 5.1: Predicción con modelo inicial para *Planta 3* Ciclo 26

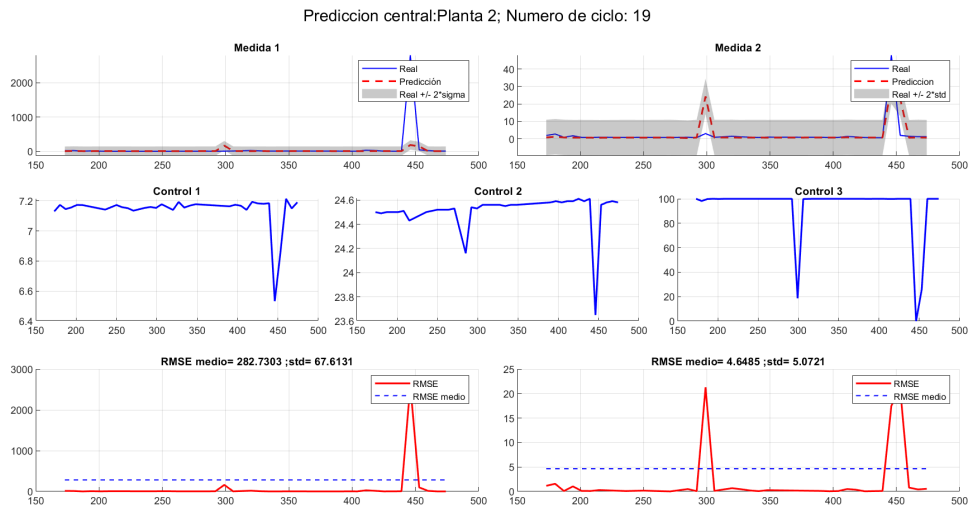


Figura 5.2: Predicción con modelo inicial para *Planta 2* Ciclo 19

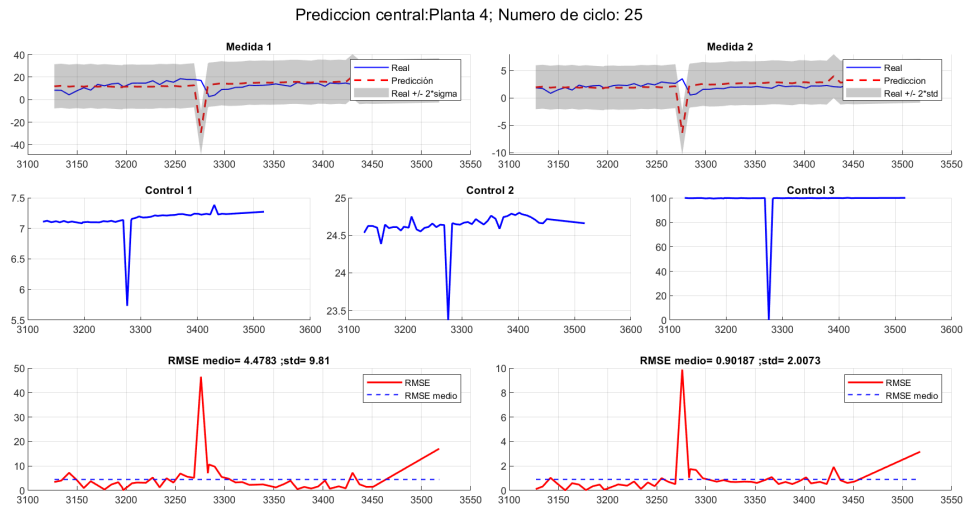


Figura 5.3: Predicción con modelo inicial para *Planta 4* Ciclo 25

Para la visualización en detalle de todos los resultados obtenidos se recomienda la lectura del Anexo D de este documento, sin embargo, su lectura no es necesaria para el correcto seguimiento del trabajo, al exponerse a continuación las conclusiones principales obtenidas tras el análisis de dichos resultados:

1. En todas las centrales podemos observar que, independientemente del valor RMSE medio obtenido para cada ciclo, la señal predicha no se ajusta a la señal real. Para las centrales *Planta 1* y *Planta 2*, las señales *Medida* presentan cambios abruptos aparentemente debidos a cambios repentinos en las variables *Control* (Véase la figura 5.2 los valores mínimos de las señales *Control* capturados en los instantes temporales cercanos a 400). Estos cambios no son predichos por el modelo con precisión, lo que supone un valor RMSE extremadamente alto. Sin embargo, para las centrales *Planta 2*, *Planta 4* y *Planta 5*, la ausencia de estos valores atípicos, que marcan el comportamiento de las señales en las centrales anteriores (ver sección 4.1.1), produce una mejor predicción que coincide en muchos casos con la tendencia de la señal real.
2. En cuanto al margen de confianza, remarcar el amplio margen de confianza de la predicción obtenidos para las centrales *Planta 1* y *Planta 2* que suponen una predicción menos precisa que el resto de plantas. Adicionalmente, debido a la simetría del margen con respecto a la media de la predicción, es habitual para todas las centrales la obtención de márgenes que incluyen valores negativos de las variables *Medida*, lo cual representa un grave error, pues esto carece de sentido dada la naturaleza química de las señales (que en general no admiten valores negativos). Las predicciones de la figura 5.3 representan un claro ejemplo de dicho comportamiento.
3. Por último, se observan múltiples casos en los que el valor real no se encuentra comprendido dentro del margen de confianza, véase las predicciones de la señal *Medida 1* de las figuras 5.2 y 5.1. Esto supone un aspecto a mejorar en la predicción.

Centrándonos en los valores medios de la medida RMSE para cada central. Se presenta la siguiente tabla:

	Planta 1	Planta 2	Planta 3	Planta 4	Planta 5
Nº de ciclos	5	6	3	4	4
RMSE Medida 1	45.372	102.219	34.751	4.650	16.337
RMSE Medida 2	2.111	2.764	4.692	1.039	2.690

Cuadro 5.1: Resultados obtenidos primera aproximación

Podemos observar que el error en la predicción de la señal *Medida 1* es considerablemente mayor que el obtenido para *Medida 2* debido a la gran diferencia entre el rango dinámico de ambas señales (ver sección 4.1.1).

Teniendo en cuenta que las centrales pueden ser divididas en dos grupos en función de los protocolos que usan, (se tiene por un lado las centrales *Planta 1* y *Planta 2*, y por otro *Planta 3*, *Planta 4* y *Planta 5*) se esperaban resultados similares entre las centrales de ambos grupos, obteniendo el primer grupo de centrales un error mayor causado por el gran número de valores atípicos en sus señales. Sin embargo, la central *Planta 3* representa un caso especial dentro del segundo grupo de centrales, debido a su alto error en las predicciones. Este error es debido a que la señal *Control 2* presenta valores completamente diferentes en un ciclo (ciclo 25) con respecto a los otros dos, lo que genera un alto error de predicción en dicho ciclo que afecta considerablemente al error medio de la central.

## 5.2. Resultados Red Bayesiana Final

Pese a no ser necesaria su lectura para el seguimiento de este trabajo, el anexo E contiene las representaciones de los resultados obtenidos mediante la utilización de la Red Bayesiana Gaussiana final descrita en la sección 3.2. Es decir, estos resultados se obtienen tras la implementación de todos los cambios propuestos en el apartado 4.3, es decir, la eliminación de la dependencia entre las variables *Medida 1* y *Medida 2*, la eliminación de los valores erróneos de *Control 1*, la utilización de la derivada de *Control 2*, la gaussianización de los datos, la inclusión de *Medida 3* y de *Control 4* al modelo, y la interpolación de las variables *Control*.

Las representaciones propuestas para el análisis son similares a las vistas en el apartado anterior, incluyendo, en la fila de las variables a predecir, la señal *Medida 3*, y en las de control, la variable *Control 4*. Adicionalmente, en las gráficas de las señales *Control*, se añaden los valores interpolados como puntos rojos dentro de la señal.

A continuación se presentan las figuras 5.4, 5.5 y 5.6, ejemplos de los resultados obtenidos y las conclusiones extraídas del análisis de los resultados.

Analizando visualmente los resultados obtenidos, se constata que, al igual que ocurría anteriormente, los valores medios de las predicciones no se ajustan bien a los valores reales de las señales predichas. A modo de ejemplo, se puede ver en la figura 5.4 como la señal predicha para todas las variables es similar al valor medio de la señal real. En el caso de las señales obtenidas en las centrales *Planta 1* y *Planta 2*, el error en las predicciones viene profundamente condicionado por los valores abruptos, no predichos, en las señales *Medida*. En la figura 5.4 se presenta un ejemplo en el cual se puede observar como el error debido al máximo atípico de todas las señales *Medida* en torno al instante  $t = 2300$  condiciona completamente el RMSE medio del ciclo, incluso en casos como el de la predicción *Medida 3*, en el que los valores máximos son predichos pero con muy poca precisión.

Este error en la predicción puede deberse a que las señales *Control* propuestas por *La Empresa*, no sean suficientes para predecir el comportamiento de las señales *Medida*, que dependerán de otras variables de la química del *primario*, no controlables o no tenidas en cuenta por *La Empresa*.

Adicionalmente, tal como se observa en las figuras 5.4, 5.5 y 5.6, el margen dinámico de la predicción, debido al proceso de gaussianización de los datos de entrada, deja de ser simétrico como en el apartado anterior, y abarca siempre valores positivos, lo que hace que este margen sí tenga un significado químico interpretable por *La Empresa*.

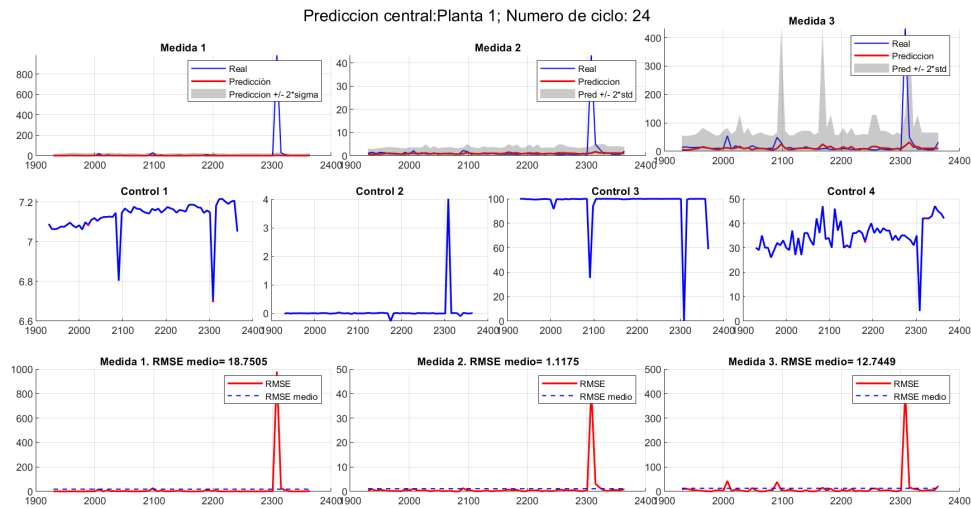


Figura 5.4: Predicción final para *Planta 1* Ciclo 24

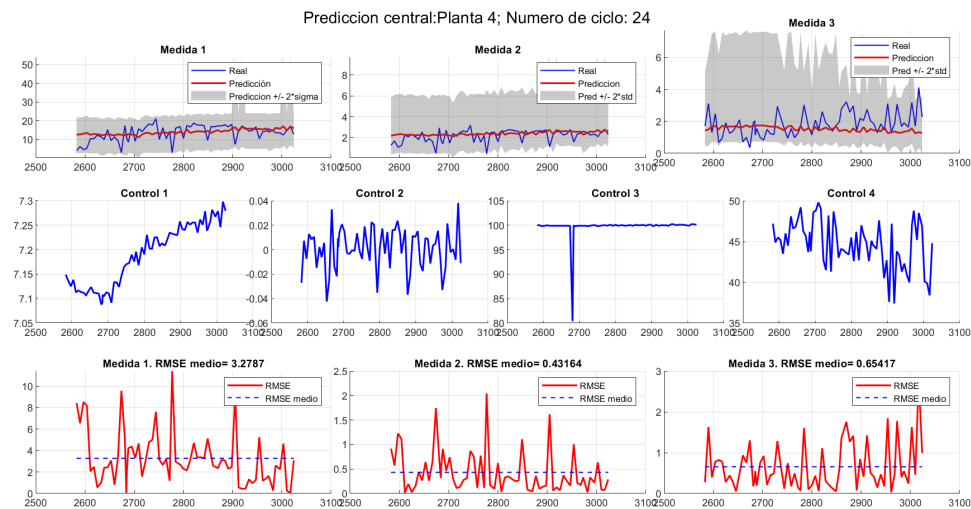


Figura 5.5: Predicción final para *Planta 3* Ciclo 25

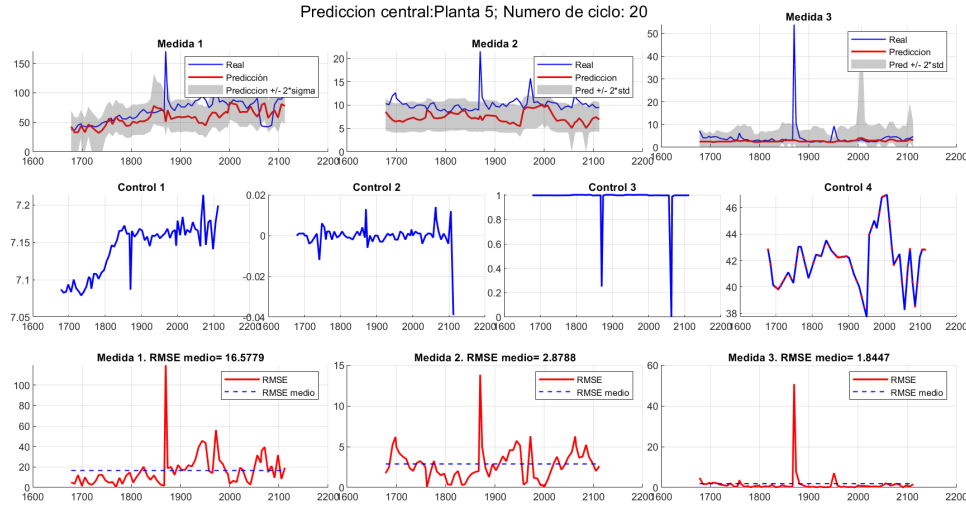


Figura 5.6: Predicción final para *Planta 5* Ciclo 20

Analizando el impacto de cada una de las modificaciones realizadas sobre la red inicial, y los datos disponibles, tal como se muestra en la figura se presenta la figura 5.7, pueden ser extraídas las siguientes conclusiones:

1. La eliminación de la dependencia entre las variables *Medida 1* y *Medida 2* no tiene ningún impacto en la predicción conjunta de ambas pues son independientes condicionalmente dadas todas las variables *Control*:

$$(M1 \perp M2)|(C1, C2, C3, C4) \rightarrow P(M1, M2|C1, C2, C3, C4) = P(M1|C1, C2, C3, C4)P(M2|C1, C2, C3, C4) \quad (5.2)$$

siendo  $M$  la correspondiente señal *Medida* y  $C$  la correspondiente señal *Control*.

2. La eliminación de los valores erróneos de *Control 1* no tiene un gran impacto en el RMSE obtenido debido al reducido número de muestras erróneas tras la inclusión de las nuevas variables al modelo, que como ha sido explicado anteriormente, reduce el número de datos de entrenamiento por la necesidad de tener datos observados de todas ellas.

El bajo impacto o nulo de ambas modificaciones sobre el RMSE medio obtenido por central es la razón por la cual no se encuentran representadas en la figura 5.7.

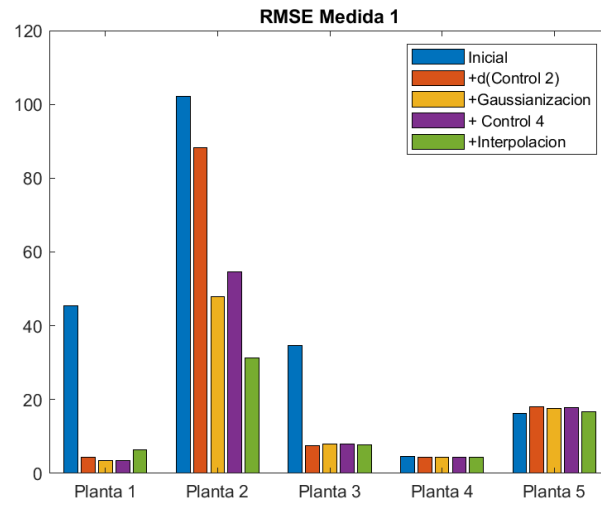
3. La primera modificación con un notable impacto en la predicción es la utilización de la señal diferenciada (o derivada discreta) de *Control 2* en lugar de su propia señal. Como se puede ver en la figura 5.7, la utilización de la derivada de *Control 2* reduce considerablemente el valor del error predicho para las centrales *Planta 1*, *Planta 2* y *Planta 3*, mientras que para las dos centrales restantes el impacto es negativo, especialmente en la predicción de la señal *Medida 2*, para la que se obtiene un incremento del 40 % del error inicial.

Cabe destacar que, para *Planta 3*, cuya predicción, como vimos anteriormente, se veía negativamente influida por la gran diferencia entre los valores de la señal *Control 2* de los diferentes ciclos, la utilización de su derivada genera la reducción de un 79 % y un 71 % del error RMSE inicial, para *Medida 1* y *Medida 2* respectivamente

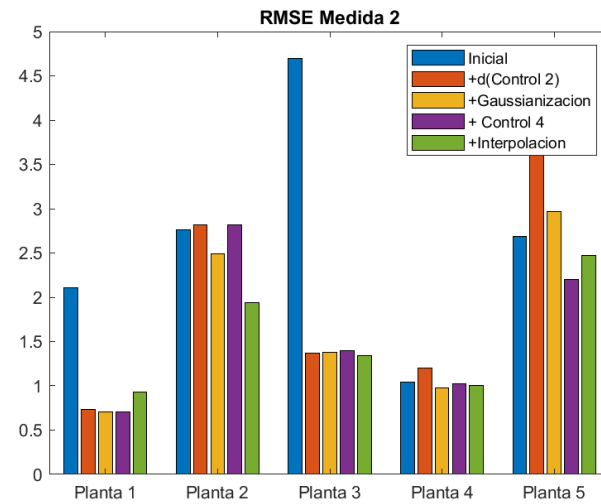


4. De la misma forma, la gaussianización de los datos de entrada a los modelos permite una mejora en todas las medidas para todas las centrales, permitiendo, en el caso de las centrales en las que la utilización de la derivada de la señal *Control 2* aumentaba el error, el nivelado del RMSE con respecto a los resultados iniciales. En concreto, esto se produce en la predicción de *Medida 2* en *Planta 4* y *Planta 5*. De esta forma, se confirma la hipótesis inicial, expuesta en 2.5, según la cual, la gaussianización de los datos permite una mejor representación de los datos disponibles a partir del modelo gráfico.
5. Con respecto al error obtenido en la predicción de la señal *Medida 3*, cabe destacar que, al haber sido incluida esta variable con posterioridad al modelo, no se ha realizado el análisis del impacto en la predicción tras las modificaciones anteriores. De esta forma, la predicción inicial realizada sobre esta variable ha sido a partir de los datos transformados y a partir de las variables *Control 1*, *Control 2* y *Control 3*. En los resultados obtenidos, se puede ver claramente la distinción entre los dos grupos de centrales ya que, el error medio obtenido para las centrales *Planta 1* y *Planta 2* es considerablemente mayor al obtenido para el resto de las centrales.
6. Como se puede observar en la figura, la inclusión de la variable *Control 4* al modelo, no supone una mejora notable para ninguna de las predicciones, exceptuando en *Planta 5*, en la cual se produce una mejora tanto en la predicción de *Medida 2* y *Medida 3*, y en la *Planta 2*, en la cual se obtiene un mayor error de predicción para todas las variables *Medida*.
7. Por último, en cuanto a la interpolación, tal y como se describe en la sección 4.3.3, se ha realizado únicamente de las variables *Control*. De esta forma se pretende evitar el entrenamiento de la red con datos erróneos causados por la interpolación errónea de los valores atípicos, propios de las señales *Medida* en determinadas centrales.

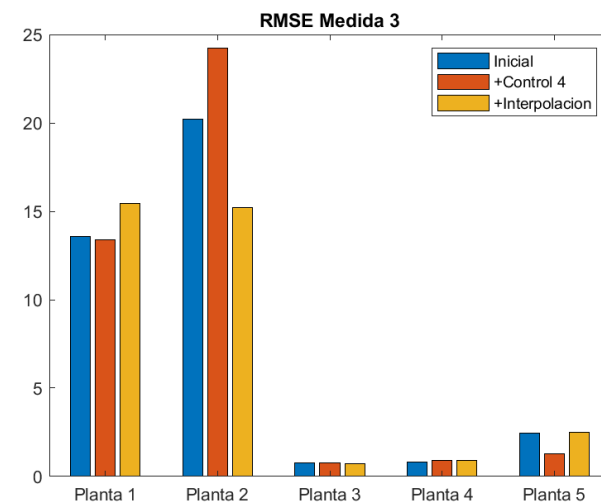
Tras el análisis del impacto de la interpolación de las variables *Control*, se concluye que este proceso produce un incremento del error de predicción en los ciclos en los cuales se realiza la interpolación de múltiples valores de la variable *Control 4*. Esto ocurre en las centrales *Planta 1* y *Planta 4*, cuyo error en la predicción de todas las señales *Medida* aumenta con respecto a la utilización de datos sin interpolar. Un ejemplo es presentado en la figura 5.6, en la cual el máximo de interpolación coincide con los instantes temporales en los que se realiza mayor interpolación de la señal *Control 2*. Por otro lado, la interpolación en la central *Planta 2* genera una reducción del error para todas las señales *Medida*. Este último hecho a esperar, pues como se vio en la sección 4.1.1, dicha señal es la más variable de todo el conjunto.



(a) Medida 1



(b) Medida 2



(c) Medida 3

Figura 5.7: Comparativa RMSE medio tras las modificaciones propuestas

### 5.3. Resultados Red Bayesiana Dinámica

Como se comentó en la sección 4.4, la aproximación dinámica propuesta en este trabajo se basa en la inclusión de las dependencias de las variables *Medida* en el instante  $t$  con respecto a todas las variables de la red en el instante  $t - N$ .

De esta forma las primeras pruebas realizadas tienen como objetivo la selección del valor de  $N$  con el que se obtienen mejores resultados en la predicción, siguiendo el mismo esquema de evaluación que el descrito en 3.3. Como se muestra en la figura 5.4, se ha probado desde la utilización del instante anterior ( $t - 1$ ) hasta el instante ( $t - 10$ ).

De este experimento, tal como refleja la figura 5.8, se puede concluir que para todas las centrales salvo para *Planta 2*, cuyo error mínimo para todas las variables *Medida* se obtiene con  $N = 7$ , el error de predicción medio aumenta con el valor de  $N$ . Es decir, los mejores resultados medios de predicción se obtienen utilizando el Modelo Bayesiano Dinámico descrito en la sección 3.3 siendo las variables latentes del modelo dependientes del estado de todas las variables en el instante previo. Y, por lo tanto, se fijará  $N = 1$  para los futuros experimentos.

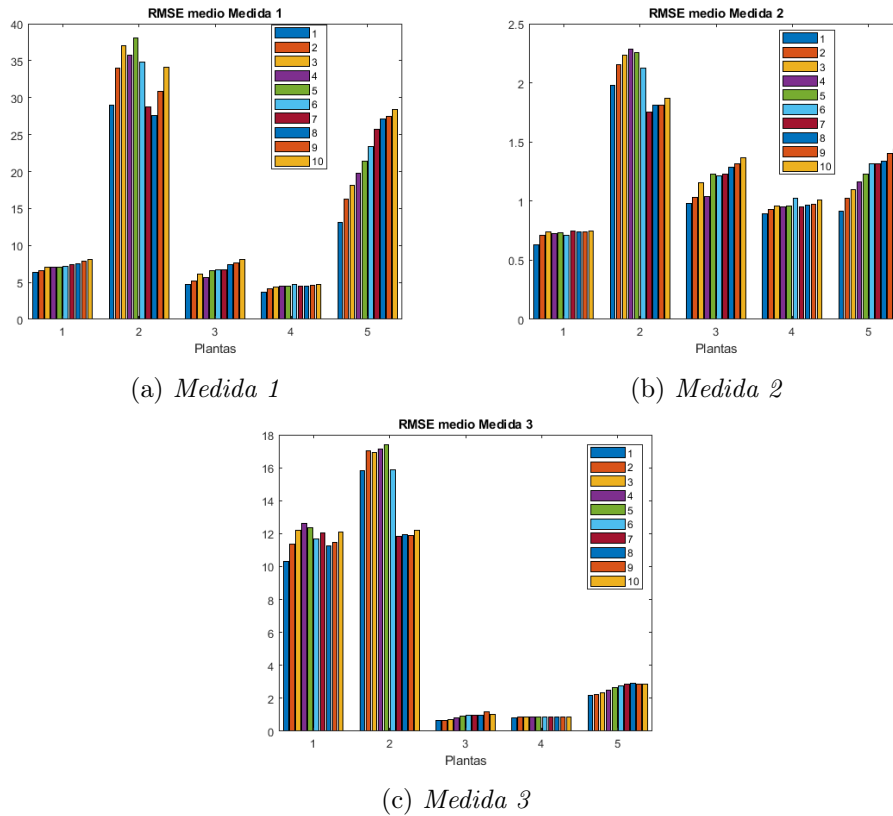


Figura 5.8: RMSE medio para todas las centrales en función de  $N$

Las figuras 5.9, 5.10 y 5.11 son ejemplos de los resultados obtenidos utilizando la Red Bayesiana Dinámica final. El resto de los resultados se encuentran en el anexo F de este documento. Se recomienda la lectura del anexo en cuestión, aunque no es necesario para el seguimiento de esta memoria.

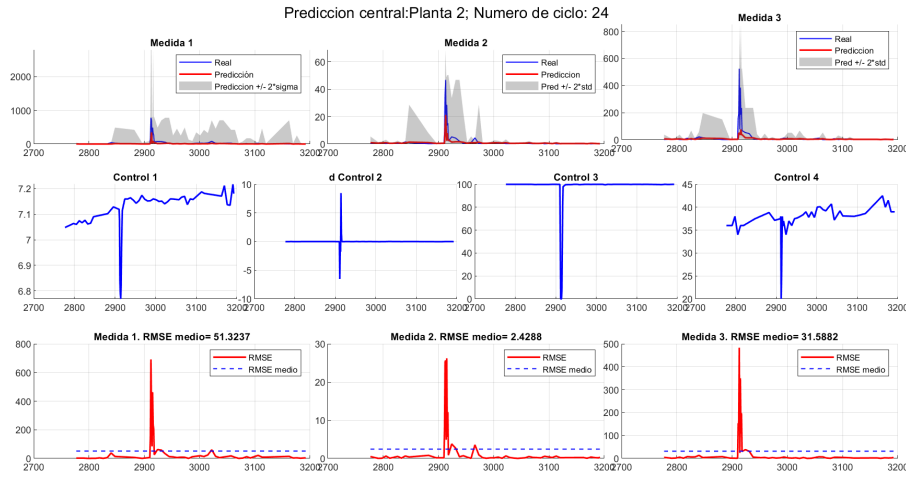


Figura 5.9: Predicción con DBN en *Planta 2* Ciclo 23

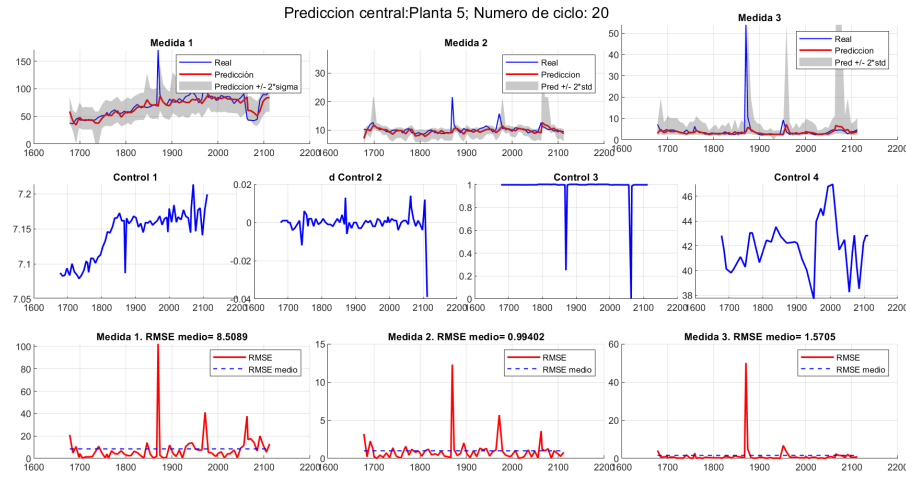


Figura 5.10: Predicción en *Planta 5* Ciclo 20

Tras el análisis visual de las predicciones, podemos observar como, al contrario que en los casos anteriores, la media de las predicciones se ajusta considerablemente a los valores reales de la señal a lo largo del ciclo. De esta forma, para las centrales *Planta 1* y *Planta 2*, generalmente, tal como se ve en la figura 5.9, el modelo es capaz de predecir valores atípicos de las variables *Medida* debidos a cambios abruptos en las variables *Control*. Algo similar ocurre, por lo general, para las variables de las centrales *Planta 3*, *Planta 4* y *Planta 5*, en las que, como podemos ver en figura 5.10, la media de la predicción se ajusta correctamente a los valores reales de la señal exceptuando valores atípicos o cambios abruptos en la señal.

Este comportamiento se estima que es debido, principalmente, a la alta correlación entre el valor de las señales *medida* en un instante, y el valor de dichas señales en el instante anterior.

Tras esto, se realiza el análisis comparativo del error de predicción del conjunto de experimentos desarrollados. Cabe destacar que, como se puede ver en la figura 5.11, debido a que tras la interpolación de los datos para el entrenamiento de la red, las centrales en *Planta 1* y *Planta 4* se obtenían peores valores de RMSE, se ha decidido realizar pruebas entrenando la red con y sin interpolación de los datos de entrada.

Como se puede ver en la figura 5.11, los mejores resultados son obtenidos utilizando la Red Bayesiana Dinámica entrenada con los datos interpolados para todas las centrales salvo para *Planta 1*, en la cual, la predicción de las variables Medida empeora para cada una de las variables.

Todavía no ha sido identificada la razón del aumento del error RMSE medio obtenido para la *Planta 1*. Pero se estima que podría estar relacionada con una mayor inestabilidad de las señales de *Control*, lo cual daría lugar a errores más relevantes al interpolar.

De esta forma, los resultados obtenidos en este apartado, representan los resultados finales obtenidos para *La Empresa*. A modo de conclusión de este capítulo, se presentan los cuadros 5.2, 5.3 y 5.4 donde se encuentra la comparativa entre los valores RMSE en la predicción obtenidos inicialmente, y los obtenidos al final de éste TFM.

	Planta 1	Planta 2	Planta 3	Planta 4	Planta 5
RMSE Inicial	45.372	102.219	34.751	4.650	16.337
RMSE Final	6.320	27.953	4.788	4.017	12.846
Diferencia (%)	86.06	72.65	86.221	13.72	21.36

Cuadro 5.2: Resumen resultados Finales *Medida 1*

	Planta 1	Planta 2	Planta 3	Planta 4	Planta 5
RMSE Inicial	2.111	2.764	4.692	1.039	2.690
RMSE Final	0.600	1.602	0.974	0.941	0.945
Diferencia (%)	71.57	42.00	79.24	9.432	64.832

Cuadro 5.3: Resumen resultados Finales *Medida 2*

	Planta 1	Planta 2	Planta 3	Planta 4	Planta 5
RMSE Inicial	13.578	20.191	0.765	0.871	2.457
RMSE Final	10.555	13.735	0.664	0.826	2.225
Diferencia (%)	22.256	31.969	13.20	5.16	9.44

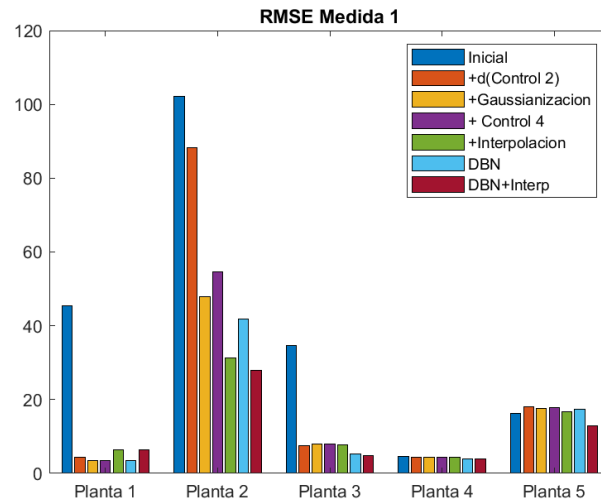
Cuadro 5.4: Resumen resultados Finales *Medida 3*

Así pues, podemos observar que, la Red Bayesiana Dinámica propuesta mejora los resultados obtenidos inicialmente para todas las medidas y todas las centrales nucleares, lo cual significa que las hipótesis iniciales del grupo, en cuanto a que la utilización de:

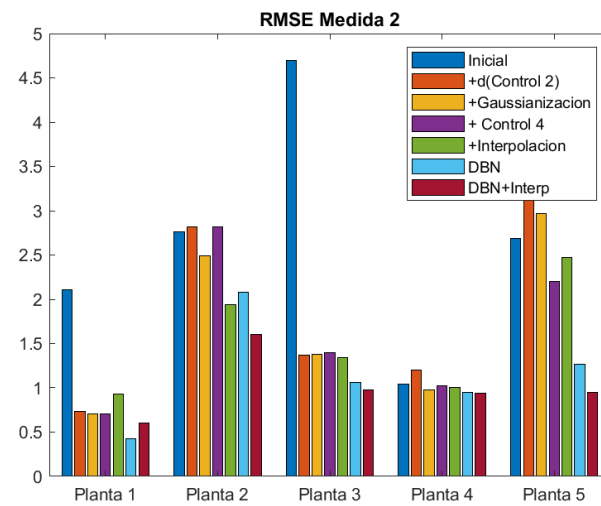
1. Técnicas de gaussianización sobre el conjunto de datos para una mejor representación, por parte del modelo gaussiano, de los mismos.
2. Interpolación de las señales *Control*, para combatir la pérdida de datos debida a las diferentes transformaciones desarrolladas.
3. Implementación de una Red Bayesiana Dinámica, para añadir dependencias con instantes temporales previos.

mejoraría los resultados iniciales, son correctos. Destacar que la centrales en las cuales se ha producido una mayor mejora del error de predicción son las centrales *Planta 1* y *Planta 3*, alcanzando, en ambos casos, una reducción del error de hasta un 86 % para la predicción de la Medida 1, y una reducción del error mayor de un 70 % para la predicción de *Medida 2*. Sin

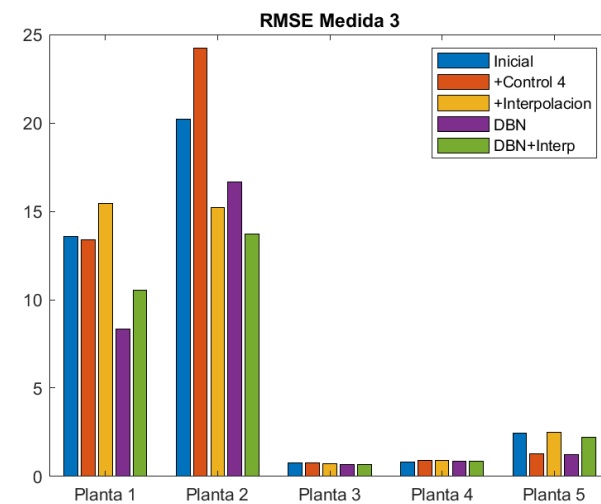
embargo, para la central *Planta 4* se ha conseguido una mejora mínima tras la implementación de todas las técnicas del desarrollo.



(a) Medida 1



(b) Medida 2



(c) Medida 3

Figura 5.11: Comparativa RMSE medio para todas las centrales para todas las técnicas





# 6

## Conclusiones y trabajo futuro

Mediante este Trabajo de Fin de Máster se ha proporcionado a *La Empresa* un análisis de las señales involucradas en la química del circuito primario del reactor nuclear que, a día de hoy, no realizaban. Adicionalmente, se ha proporcionado un Software, que, mediante la utilización de Redes Bayesianas Gaussianas y Redes Bayesianas Dinámicas Gaussianas, proporciona mayor conocimiento acerca de las relaciones entre dichas variables, además de permitir la predicción de las señales de interés a partir de las señales controlables, para cada una de las centrales de cuya gestión se encarga *La Empresa*. De esta forma, se cumplen los objetivos iniciales del proyecto.

Adicionalmente, con el objetivo de mejorar el rendimiento de la Red Bayesiana Gaussiana inicial, se han realizado:

1. Transformaciones, propuestas por *La Empresa*, sobre el conjunto de datos de entrenamiento.
2. Gaussianización de los datos de entrada mediante la ecualización de histogramas, para un mejor representación de los datos mediante el modelo gaussiano.
3. Interpolación lineal de los datos, para solventar, de esta forma, la pérdida de datos debida a las diferentes transformaciones.
4. Implementación de una Red Bayesiana Dinámica basada en la Red Bayesiana inicial, para incluir dependencias con instantes temporales previos.

De esta forma, mediante la integración de todos los cambios desarrollados, se alcanzan los mejores resultados del proyecto, que suponen una mejora considerable en la predicción con respecto al modelo desarrollado inicialmente.

Pese a que los resultados finales satisfacen las expectativas iniciales del proyecto, se considera que hay un margen de mejora alcanzable mediante la realización de diferentes tareas en el futuro, como por ejemplo:

- Utilización de técnicas más complejas de gaussianización de datos, por ejemplo *Normalizing Flows*.

- Estudio de técnicas de aprendizaje de la estructura de la red que permitan la generación de un modelo más complejo que incluya todas las variables que influyan a las variables de interés.
- Incremento de las dependencias temporales de la Red Bayesiana Dinámica: Aumentar la complejidad de la red, incluyendo la dependencia con más instantes temporales.
- Generación de la interfaz gráfica final, añadiendo mayor funcionalidad.
- Adicionalmente, para ampliar el conocimiento de *La Empresa* en cuanto a la generación de residuo final generado en cada ciclo en el interior del reactor, se plantea la predicción del residuo como un problema de regresión en función de métricas de las variables de interés predicha.

# Bibliografía

- [1] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [2] William M Bolstad and James M Curran. *Introduction to Bayesian statistics*. John Wiley & Sons, 2016.
- [3] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [4] David Barber. *Bayesian reasoning and machine learning*. Cambridge University Press, 2012.
- [5] Michael I Jordan et al. Graphical models. *Statistical Science*, 2004.
- [6] Morris L Eaton. Multivariate statistics: a vector space approach. *Jhon Wiley & Sons, Inc*, 1983.
- [7] Theodore Wilbur Anderson and Ingram Olkin. Maximum-likelihood estimation of the parameters of a multivariate normal distribution. *Linear algebra and its applications*, 1985.
- [8] In Jae Myung. Tutorial on maximum likelihood estimation. *Journal of mathematical Psychology*, 2003.
- [9] Kevin Patrick Murphy and Stuart Russell. Dynamic bayesian networks: representation, inference and learning. 2002.
- [10] Zoubin Ghahramani. Learning dynamic bayesian networks. In *International School on Neural Networks, Initiated by IIASS and EMFCSC*. Springer, 1997.
- [11] V Mihajlovic and Milan Petkovic. Dynamic bayesian networks: A state of the art. *University of Twente Document Repository*, 2001.
- [12] Scott Saobing Chen and Ramesh A Gopinath. Gaussianization. In *Advances in neural information processing systems*, pages 423–429, 2001.
- [13] George Saon, Satya Dharanipragada, and Daniel Povey. Feature space gaussianization. In *2004 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pages I–329. IEEE, 2004.
- [14] Jason Pelecanos and Sridha Sridharan. Feature warping for robust speaker verification. 2001.
- [15] Kevin Murphy et al. The bayes net toolbox for matlab. *Computing science and statistics*, 2001.
- [16] Myles Hollander, Douglas A Wolfe, and Eric Chicken. *Nonparametric statistical methods*, volume 751. John Wiley & Sons, 2013.
- [17] Jean Dickinson Gibbons and Subhabrata Chakraborti. *Nonparametric statistical inference*. Springer, 2011.